# The Impact of Delay on the Diversity, Multiplexing, and ARQ Tradeoff

Tim Holliday
Princeton University

Andrea Goldsmith
Stanford University

H. Vincent Poor
Princeton University

*Abstract*— A substantial amount of research has focused on analyzing and achieving the diversity-multiplexing tradeoff in multiple antenna (MIMO) wireless communications. Recently, ARQ protocols have been added to these formulations and shown to perform as a type of diversity. Our goal in this paper is to find the optimal operating point in the diversity-multiplexing-ARQ tradeoff, with a particular focus on delay sensitive systems. Previous results in this area construct performance measures through the use of high SNR asymptotic approximations. While effective, these approximations tend to trivialize the delay performance of MIMO systems. We present a dynamic programming formulation for finding the optimal diversity gain, multiplexing gain, and ARQ window size, without relying on a high SNR approximation. Our results show that the a delay sensitive system requires one to adapt diversity and multiplexing to the time-varying workload in the system. We provide numerical examples that demonstrate the significant performance gains that can be achieved by choosing an adaptive policy over a static allocation of diversity and multiplexing.

## I. INTRODUCTION

Multiple antennas can significantly improve the performance of wireless systems. Roughly speaking, signalling schemes that exploit the spatial domain can be used to increase the data rate of the system (through spatial multiplexing) or to decrease the probability of error (through spatial diversity). The work by Zheng and Tse [13] demonstrated that both diversity and multiplexing can be accomplished simultaneously. However, there is a fundamental tradeoff between the two quantities: higher spatial multiplexing gain leads to lower diversity and vice versa. Recently, El Gamal and Caire [2] extended the diversity-multiplexing tradeoff to include ARQ protocols. In ARQ, the receiver feeds back a message to the transmitter denoting whether or not the transmission was successful. If the message was received in error, the transmitter re-sends another version of the message until it is received correctly. The results in [2] show that the ARQ process can be viewed as another form of diversity. Moreover, they show that a sufficiently large ARQ window permits the recovery of nearly all the spatial diversity lost to spatial multiplexing. Hence, through the use of ARQ one can "flatten out" the diversity-multiplexing tradeoff curve and develop systems with essentially both full multiplexing and full diversity.

The results in [2] provide a number of interesting suggestions for system design. Perhaps the most surprising is that large ARQ windows provide substantial diversity gains without destroying multiplexing gains (i.e. ARQ retransmissions do not cause a significant reduction in throughput performance). Hence, if the performance metric of interest is link layer throughput then the optimal design choice is to use a large ARQ window and as much multiplexing as possible.

In [9] we considered the question: "Given the diversity-multiplexing tradeoff region, where should one choose to operate?" We showed that a standard end-to-end distortion metric allowed one to rigorously determine the optimal level of multiplexing and diversity to be used in the design of a multiple antenna system. The new results from [2] also apply to the distortion problem in [9]. One can easily show that the optimal average distortion is achieved by choosing a high multiplexing gain and utilizing all of the available ARQ window to create diversity.

Our goal in this paper is to consider a different distortion measure that also accounts for end-to-end delay. This performance metric is more appropriate for voice, video, and other types of multimedia traffic. We will show that a delay sensitive distortion measure results in a fundamentally different type of optimal control in the diversity, multiplexing, and ARQ tradeoff; where high levels of multiplexing and long ARQ windows are no longer optimal. Indeed, we show that the optimal amount of multiplexing and ARQ window length are both highly dependent on the traffic's delay sensitivity. We also show that substantial gains can be realized by adapting the choice of diversity and multiplexing to the time-varying system workload.

In the next section we review previous results on the diversity-multiplexing-ARQ tradeoff. We examine the impact of ARQ on this tradeoff and consider some of the problems these formulations introduce when we wish to examine delay. Specifically, we show that the high SNR asymptotic regime can trivialize many measures of delay. In Section III we introduce a distortion model that also accounts for end-to-end delay. We model an encoder concatenated with a MIMO-ARQ channel as a large Markov chain. Then we show that we can formulate and solve a Markov decision process to minimize this new distortion measure and find the optimal adaptive diversity-multiplexing control. Section IV contains some initial numerical examples. We conclude in Section V.

## II. PREVIOUS WORK REGARDING THE DIVERSITY-MULTIPLEXING-ARQ TRADEOFF

We will use the same channel model and notation from [2]. Consider a wireless link with $M$ transmit antennas and $N$ receive antennas. We assume that the system performs the following ARQ scheme. Each information message is encoded into a sequence of $L$ blocks each of size $T$. Transmission

commences with the first block and after decoding the message the receiver sends a positive (ACK) or negative (NACK) acknowledgement back to the transmitter. In the case of a NACK the transmitter sends the next block in the sequence and the receiver uses all accumulated blocks to try to decode the message. This process proceeds until either the receiver correctly decodes the message or until all $L$ blocks have been sent. If a NACK is sent after the transmission of the $L$th block then an error is declared, the message is removed from the system, and the transmitter starts over with the next queued message. For the sake of clarity we will use the term "block" to describe a single transmission attempt during the ARQ process. We will refer to an attempt to send a message using ARQ as an "ARQ round". Hence, each round of ARQ consists of up to $L$ blocks, and each block is of size $T$.

The fading coefficients $h_{ij}$ that model the gain from transmit antenna $i$ to receive antenna $j$ are i.i.d. complex Gaussian with unit variance. The channel gain matrix $H$ with elements $H(i,j) = (h_{ij} : i \in \{1, \dots N\}, j \in \{1, \dots, M\})$ is assumed to be known at the receiver and unknown at the transmitter. We assume that the channel remains constant over the entire ARQ round of up to $LT$ symbols, while each ARQ round is i.i.d. Therefore, in block $l \in \{1, \dots, L\}$ of an ARQ round we can represent the channel as

$$Y_l = \sqrt{\frac{SNR}{M}} H X_l + W, \tag{1}$$

where $X_l \in \mathcal{C}^{M \times T}$ and $Y_l \in \mathcal{C}^{N \times T}$ are the transmitted and received signals in block $l$, respectively. The additive noise vector $W$ is i.i.d. complex Gaussian with unit variance.

### A. The Fundamental Tradeoff

With the above model in hand define a family of codes $\{C(SNR)\}$, indexed by the $SNR$ level. Each code has length $LT$ and the bit rate of the first block in each code is $b(SNR)/T$. Suppose we consider a sequence of ARQ rounds. At time $s$ the random variable $B[s] = b(SNR)$ if a message is successfully decoded at the receiver, and $B[s] = 0$ otherwise. Then, we can define the average throughput of the ARQ protocol using these codes as

$$\eta(SNR) = \liminf_{\tau \to \infty} \frac{1}{T\tau} \sum_{s=1}^{\tau} B[s], \tag{2}$$

and we can view $\eta(SNR)$ as the average number of transmitted bits per channel use. Further define $P_e(SNR)$ as the average probability of error of the ARQ round (i.e. the probability that a NACK is sent after $L$ transmission rounds). Define the multiplexing gain of the ARQ protocol as

$$r = \lim_{SNR \to \infty} \frac{\eta(SNR)}{\log SNR}, \tag{3}$$

and the diversity gain as

$$d = -\lim_{SNR \to \infty} \frac{\log P_e(SNR)}{\log SNR}. \tag{4}$$

For each $r$ and $L$ we define the optimal diversity gain $d^*(r, L)$ as the supremum of the diversity gain achieved by any scheme. For $L = 1$ (i.e. no ARQ) we have the following result from [13].

**Diversity-Multiplexing Tradeoff:** Assume the block length $T \geq M + N - 1$. Then the optimal tradeoff between diversity gain and multiplexing gain is given by $d^*(r, 1) = f(r)$, where $f(r)$ is the piecewise linear function joining the points $(k, (M - k)(N - k))$, at integer values of $k$ for $0 \leq k \leq \min(M, N)$.

For $L > 1$ we have the following result from [2].

**Diversity Gain of ARQ:** The diversity gain for the ARQ protocol with a maximum of $L$ blocks is

$$d^*(r, L) = f\left(\frac{r}{L}\right). \tag{5}$$

The diversity gain achieved by ARQ is quite remarkable. According to (5), for any $r < \min(M, N)$ we can achieve the full diversity gain $d = MN$ for sufficiently large $L$. This suggests that there is little point in utilizing spatial diversity since we can always acquire any needed diversity through ARQ.

In order to analyze the diversity, multiplexing, and ARQ tradeoff in delay sensitive systems we must recognize two important conditions of the above results. First, in delay sensitive systems we may not be able to tolerate a long ARQ window (in some cases ARQ may not be tolerated at all). Second, we must carefully consider the impact of the high SNR asymptotic regime, which is crucial in the proofs of the above results. Specifically, in the high SNR regime the occurrence of a NACK in the ARQ protocol becomes a rare event (i.e. the probability of a NACK tends to zero as $SNR \to \infty$). Therefore, with probability tending to one, each message is decoded correctly during the first transmission attempt – resulting in a multiplexing gain equivalent to that of a system without ARQ. The increasingly rare errors are corrected by the ARQ process, which results in increased diversity.

The main difficulty in using these asymptotic results to evaluate delay performance is that in the high SNR regime there is essentially *no delay* . Using standard results from queueing theory one can show that the rare error events will cause arriving messages to almost always find the system empty. Hence, with high probability an arriving message will immediately begin transmission and suffer no queueing delay. In wireless systems, errors during a transmission attempt are not rare events. Indeed, modern wireless systems typically become reliable only after the application of ARQ. In other words, errors after completion of the ARQ process might be rare events, but errors *during* the ARQ process are not rare. As we shall see in the next section, this subtle difference results in an entirely different optimization problem and solution when we consider delay sensitive distortion measures.

## III. A Delay-Distortion Model for the Diversity-Multiplexing-ARQ Tradeoff

This section presents our system model for minimizing a delay sensitive end-to-end distortion measure. We do not assume a high SNR regime in this analysis. However, we do assume that the $SNR$ is fixed for each problem instance (i.e. we do not optimize power control).

We assume the original source data $u$ is a random vector with probability density $f(u)$, which has support on a closed, bounded subset of $\Re^k$ with non-empty interior. During each transmission block of length $T$ an instance of $u$ arrives at the system independently with probability $\lambda$ and is queued for transmission. We assume that each message has a playout deadline $k$ at the receiver. Hence, if a message arrives at time $t$ and is not received by time $t + kT$ then its deadline expires and the message is dropped from the system. We assume that each message is quantized according to the scheme discussed below. The quantized version of each message is then mapped into a codeword in the codebook $\{C(SNR)\}$ and passed to the MIMO-ARQ transmitter discussed in the previous section.

Due to the random message arrival times and the random completion times of the ARQ process we will have queueing and delay in this system. Our goal is to select a diversity gain, multiplexing gain, and ARQ window size to minimize the distortion created by both the quantizer and the messages lost due to delay. The intuition behind the diversity-multiplexing-ARQ tradeoff is straightforward. We would we like to use as much multiplexing as possible since this will allow us to use more bits to describe a message and reduce encoder distortion. However, high levels of multiplexing induce more errors in the wireless channel, thereby requiring longer ARQ windows to reduce errors. The longer ARQ windows induce higher delays, which also cause higher distortion due to messages missing their deadlines. We must balance all of these quantities to optimize system performance.

### A. Message Distortion Model

An $s$-bit quantizer is applied to each message via the following transformation:

$$F(u) = \sum_{i=1}^{2^s} v_i I_{[A_i]}, \qquad (6)$$

where $I_{[A_i]}$ is the standard indicator function, and $\{A_i\}_{i=1}^{2^s}$ is a partition of $\Re^k$ into disjoint regions. Each region $A_i$ is represented by a single codevector $v_i$.

We assume the encoder/decoder pair achieves the noiseless distortion [5]

$$D_s(F) = 2^{-ps/k}, \qquad (7)$$

where

$$D_s(F) = \sum_{i=1}^{2^s} \int_{A_i} ||u - v_i||^p f(u)du, \qquad (8)$$

and $||u - v_i||^p$ is the $p$th power of the Euclidian norm. We use this distortion measure strictly for the sake of simplicity. In-

deed, we only require a function that is convex and decreasing in $s$ and discuss this issue in more detail in [10].

Assume that the rate of the channel codebook $\mathcal{C}\{SNR\}$ is matched to the rate of the quantizer. Hence, the number of bits used to describe the message will determine $b(SNR)$ and the amount of spatial multiplexing gain in the system. Furthermore, we assume that we can achieve the optimal point on the diversity-multiplexing tradeoff curve. Therefore a choice of spatial multiplexing gain also specifies the spatial diversity gain.

We assume, as in [8], [9], [6] that the total distortion $D_T(F, SNR)$ can be split into two dependent pieces

$$D_T(F, SNR) = D_s(F) + D_e(d, SNR), \qquad (9)$$

where $D_e(d, SNR)$ is the distortion caused by messages declared in error. Here the errors are incurred whenever the ARQ process fails or when a message's deadline expires. As in [6], [8] we also assume the distortion due to erroneous messages is bounded by the overall loss probability:

$$D_e(d, SNR) \le P_e(SNR) + (1 - P_e(SNR))P\{Delay > k\}, \qquad (10)$$

where $P\{Delay > k\}$ is the probability that a message violates its deadline. We note that this assumption coincides with tight rigorous bounds [8] as well as heuristic distortion measures that have been fitted to real world traffic streams [12].

### B. Minimizing Total Distortion

Our goal is to minimize the total delay-distortion bound

$$
\begin{aligned}
D_T(F, SNR) &\le D_s(F) + P_e(SNR) & (11) \\
&+ (1 - P_e(SNR))P\{Delay > k\}.
\end{aligned}
$$

In order to optimize (11) we require a formulation that accounts for the different delays experienced by each message. Hence we turn to the theory of Markov decision processes to model and solve this problem.

Without loss of generality assume that the queue described in the previous section is of maximum size $k$. Note that each message requires at least one time block of size $T$ for transmission, hence any arriving message that sees more than $k$ messages in the queue will not be able to meet its deadline. Unlike standard queueing models that only track the number of messages awaiting transmission, we must also track the amount of time a particular message has waited. For example, if one message is queued for transmission we will differentiate between the states where a message has just arrived and a message whose deadline is about to expire. Since we may only have a finite number of messages in the queue, this combined message and waiting time model exists in a finite space. Define the queue process $X_Q = (X_Q(n) : n \ge 0)$, which takes values on a finite space $\mathcal{X}_Q$. Likewise, we may define the state of the ARQ process $X_L = (X_L(n) : n \ge 0)$ on a finite space $\mathcal{X}_\mathcal{L}$. Here, the state of the ARQ process denotes the number of the current transmission block in the current ARQ round. Define the overall state of the system as a process

1447

$X = (X(n) : n \geq 0)$ such that $X(n) = (X_Q(n), X_L(n))$ (i.e. the space $\mathcal{X}$ is the product space of $\mathcal{X}_Q$ and $\mathcal{X}_{\mathcal{L}}$).

Since the arrival process is geometric and each ARQ round is assumed to be i.i.d., the process $X$ is a finite-state discrete-time Markov chain. The transition dynamics of this Markov chain are governed by the choices of diversity, multiplexing, and ARQ window size. Assume that at the start of each ARQ round the transmitter chooses the number of bits to assign to the vector encoder and hence the amount of spatial diversity and multiplexing in the codeword selected from $\{C(SNR)\}$. The transmitter also selects the length of the ARQ window. These choices then remain fixed until either the message is received or the ARQ window expires. Define the space of actions $\mathcal{A}$ as the set of all possible combinations of multiplexing gain and ARQ window length. (Note that a choice of multiplexing gain implicitly selects the number of bits given to the source encoder as well as the amount of spatial diversity available.) We assume that the number of antennas $M$ and $N$ are finite and that the ARQ window size is also finite. Hence, the action space $\mathcal{A}$ consists of a finite set.

Define a control policy $g$ as a probability distribution on the space $\mathcal{X}$ x $\mathcal{A}$. We can view the elements of $g$ as

$$g(x, a) = P\{\text{action } a \text{ chosen in state } x\}, \ \forall x \in \mathcal{X}, a \in \mathcal{A}.$$

For any control $g$, the Markov chain $X$ is irreducible and aperiodic[1]. Define $Q(g)$ as the transition matrix for $X$ corresponding to control policy $g$. Hence, $Q(g) = (Q_{i,j}(g) : i, j \in \mathcal{X})$ is a stochastic matrix with entries

$$\begin{aligned} Q_{i,j}(g) &= P\left(X(n+1) = j | X(n) = i, g\right) \\ &= \sum_{a \in \mathcal{A}} P\left(X(n+1) = j | X(n) = i, A(n) = a\right) g(i, a)). \end{aligned}$$

For each state-action pair we define a reward function $r(x, a)$. For the states in $\mathcal{X}$ corresponding to completion of the ARQ process the reward function denotes the distortion incurred in that particular state. Hence,

$$r(x, a) = 2^{-ps/k} + I_{[\text{ARQ Fail}]} + I_{[\text{ARQ Succeed}]} I_{[Delay > k]}. \tag{12}$$

Let $\mathcal{G}$ be the set of all available control policies. Then for any $g \in \mathcal{G}$ define the limiting average value of $g$ starting from state $x$ as

$$V(x, g) = \limsup_{n \to \infty} \left[ \left(\frac{1}{n+1}\right) \sum_{k=0}^{n} E_{x,g}\left[r(X(k), g)\right] \right].$$

where $r(X(k), g)$ is the random reward earned at time $k$ under control policy $g$. Since $X$ is an irreducible and aperiodic Markov chain for any control $g$ we know [1] that the above value function reduces to

$$V(x, g) = \pi(g) r(g) \ \forall x \in \mathcal{X}, \tag{13}$$

where $\pi(g) = \pi(g) Q(g)$ is the stationary distribution of $X$ under control $g$ and $r(g)$ is the column vector of rewards

[1]To create a non-irreducible Markov chain we would require the ability to successfully transmit a packet with probability one.

earned for each state $x \in \mathcal{X}$ under control $g$. Hence, the value function is simply the expected value of our reward function $r$ with respect to the stationary distribution of $X$. Notice that given our definition for $r$ in (12), the value function $V(g)$ provides us with the delay-based distortion (11) caused by control policy $g$.

Our goal is to find a $g \in \mathcal{G}$ that minimizes $V(x, g)$. From [1] we know this problem can be solved through the following linear program.

$$\min_s \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} r(x, a) s_{xa} \tag{14}$$

subject to:

$$\sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \left(\delta(x, x') - p(x'|x, a)\right) s_{xa} = 0, \ x' \in \mathcal{X}$$

$$\sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} s_{xa} = 1,$$

$$s_{xa} \geq 0; \ a \in \mathcal{A}, \ x \in \mathcal{X},$$

where $\delta(x, x')$ is the Kronecker delta, $s_{xa}$ is the steady-state probability of being in state $x$ and taking action $a$, and $p(x'|x, a)$ is the probability of jumping to state $x'$ given action $a$ in state $x$. The state-action frequencies $s_{xa}$ provide a unique mapping to an optimal control $g^*$ [1].

With this dynamic programming formulation in hand we can solve for the optimal diversity gain, multiplexing gain, and ARQ window size as a function of queue state and deadline sensitivity. We demonstrate the performance of these solutions with a numerical example in the next section.

## IV. Numerical Example

Consider the system setup described above with messages arriving in each time block with probability $\lambda = .9$. We assume a 4x4 MIMO-ARQ system ($M = N = 4$) utilizing the incremental redundancy codes proposed in [3] which have been shown to achieve the diversity-multiplexing-ARQ tradeoff. The ARQ window size will take values in a finite set $L \in \{1, \ldots, 4\}$. We will allow the deadline parameter $k$ to range over several values ($k \in \{2, \ldots, 16\}$) in order to examine the impact of delay sensitivity on the solution to our dynamic program (14). For each value of $k$ we solve a new version of (14), the plots below contain the data accumulated from all of the solutions.

Figure 1 plots the optimal ARQ window length as a function of queue state for different values of $k$. We can see that for short deadlines we cannot afford long ARQ windows for any queue state. As the deadlines become more relaxed we can increase the ARQ window size. Although as the queue fills up we are forced to again decrease the amount of ARQ diversity.

Figure 2 plots the optimal multiplexing gain $r$ as a function of queue state for different values of $k$. Here we can see that with short deadlines we must utilize fairly low amounts of spatial multiplexing (i.e. high spatial diversity), since we cannot use ARQ diversity. As the deadlines become more relaxed we can increase the amount of spatial multiplexing
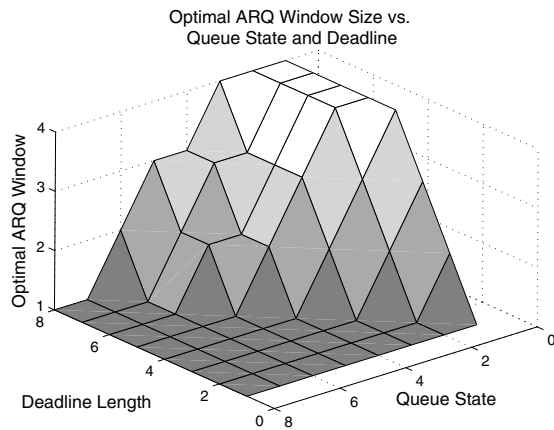
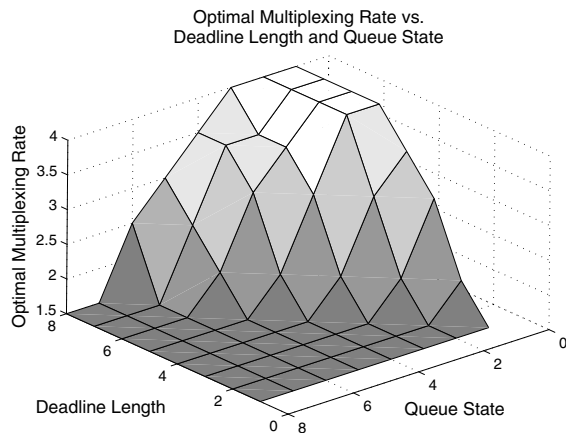Fig. 1. Optimal ARQ window size vs. queue state vs. deadline parameter $k$.



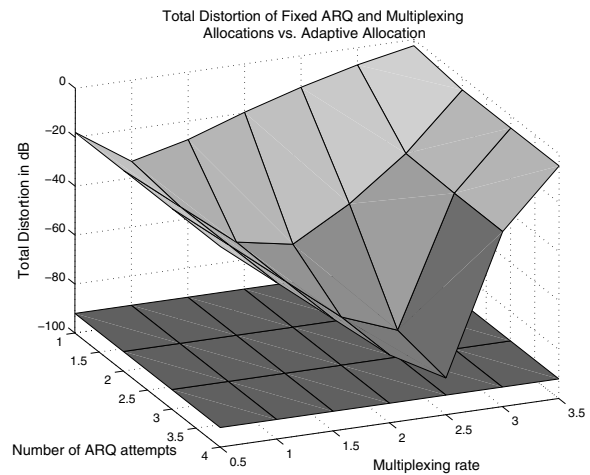Fig. 2. Optimal multiplexing gain vs. queue state vs. deadline parameter $k$.



Fig. 3. Distortion for the fixed allocation problem vs. multiplexing gain vs. ARQ window size.

be suitable for measuring delay performance. We then present a dynamic programming framework for optimizing the choices of diversity, multiplexing, and ARQ window length in delay constrained environments. Our results show that it is optimal to adapt ARQ and the diversity-multiplexing gains to the amount of data queued for transmission. Furthermore, the gains in performance achieved by adaptation can be significant when compared with static allocation policies.

and utilize ARQ for diversity. Once again, as the queue fills up we must switch back to low levels of multiplexing in order to ensure that traffic is cleared from the system on time.

We also evaluate the performance advantage gained by adapting the settings of diversity, multiplexing, and ARQ rather than choosing fixed allocations. For $k = 4$ we computed the distortion resulting from all possible fixed allocations of ARQ window length and multiplexing gain. The curved surface in Figure 3 plots the distortion of these fixed allocations for all values of $L$ and $r$. The flat surface in Figure 3 is the distortion achieved by the adaptive scheme (plotted as a reference). Even in the most favorable cases, the adaptive scheme outperforms any fixed scheme by more than 50%.

## V. CONCLUSION

In this paper we investigate the optimal diversity-multiplexing-ARQ tradeoff in terms of minimizing a delay sensitive end-to-end distortion measure. We show that current analysis, which relies on a high SNR approximation, may not

## REFERENCES

[1] Bertsekas, D. P., "Dynamic Programming and Optimal Control", *Athena Scientific* , 1995.
[2] H. El Gamal, G. Caire, and M. O. Damen, "The MIMO ARQ channel: Diversity-multiplexing-delay tradeoff", *submitted to the IEEE Trans. on Inform. Theory*, Nov. 2004.
[3] H. El Gamal, G. Caire, and M. O. Damen, "Lattice coding and decoding achieve the optimal diversity-vs-multiplexing tradeoff of MIMO channels" , *IEEE Transactions on Information Theory*, June 2004.
[4] H. El Gamal, G. Caire, and M. O. Damen, "On the optimality of incremental redundancy LAST coding", *Proc. of WCNC 2005"*, Vol. 2, pp. 1472 - 1476, June 2005.
[5] A. Gersho, "Asymptotically Optimal block quantization", *IEEE Trans. Inform. Theory* , Vol. 25, pp. 373-380, July 1979.
[6] B. Girod, K. Stuhlmuller, N. Farber, "Trade-Off Between Source and Channel Coding for Video Transmission", *Proc. of ICIP'2000*, Vol. 1, pp. 399-402, Sept. 2000.
[7] D. Gross, C.M. Harris, "Fundamentals of Queueing Theory", Wiley, New York, 2000.
[8] B. Hochwald, K. Zeger, "Tradeoff Between Source and Channel Coding", *IEEE Trans. Inform. Theory* , Vol. 43 , pp. 1412 - 1424, Sept. 1997.
[9] T. Holliday, A. Goldsmith, Joint-Source Channel coding in MIMO Systems, *Proc. of the Allerton 2004 Conf.* , Sept. 2004.
[10] T. Holliday, V. Poor, A. Goldsmith, Cross-Layer Design of MIMO Systems, *journal version in preparation* .
[11] M. Kuhn, I. Hammerstroem, A. Wittneben, "Linear Scalable Space-Time Codes: Tradeoff Between Spatial Multiplexing and Transmit Diversity", *Proc. of SPAWC 2003*, June 2003.
[12] K. Stuhlmuller, N. Farber, M. Link, and B. Girod, "Analysis of Video Transmission over Lossy Channels', *IEEE Journal on Selected Areas in communications*, vol. 18, pp. 1012-1032, June 2000.
[13] L. Zheng, D. Tse, "Diversity and Multiplexing: the Optimal Tradeoff in Multiple Antenna Channels", *IEEE Trans. Inform. Theory* , Vol. 49, pp. 1073 - 1096, May 2003.