

# Performance of Hybrid-ARQ in Block-Fading Channels: A Fixed Outage Probability Analysis

Peng Wu and Nihar Jindal

## Abstract

This paper studies the performance of *hybrid*-ARQ (automatic repeat request) in Rayleigh block-fading channels. The long-term average transmitted rate is analyzed in a fast-fading scenario where the transmitter only has knowledge of channel statistics, and, consistent with contemporary wireless systems, rate adaptation is performed such that a target outage probability (after a maximum number of H-ARQ rounds) is maintained. H-ARQ allows for early termination once decoding is possible, and thus is a coarse, and implicit, mechanism for rate adaptation to the instantaneous channel quality. Although the rate with H-ARQ is not as large as the ergodic capacity, which is achievable with rate adaptation to the instantaneous channel conditions, even a few rounds of H-ARQ make the gap to ergodic capacity reasonably small for operating points of interest. Furthermore, the rate with H-ARQ provides a significant advantage compared to systems that do not use H-ARQ and only adapt rate based on the channel statistics.

## I. INTRODUCTION

ARQ (automatic repeat request) is an extremely powerful type of feedback-based communication that is extensively used at different layers of the network stack. The basic ARQ strategy adheres to the pattern of transmission followed by feedback of an ACK/NACK to indicate successful/unsuccessful decoding. If simple ARQ or *hybrid*-ARQ (H-ARQ) with Chase combining (CC) [1] is used, a NACK leads to retransmission of the same packet in the second ARQ round. If H-ARQ with incremental redundancy (IR) is used, the second transmission is not the same as the first and instead contains some “new” information regarding the message (e.g., additional parity bits). After the second round the receiver again attempts to decode, based upon the second ARQ round alone (simple ARQ) or upon both ARQ rounds (H-ARQ, either CC or IR). The transmitter moves on to the next message when the receiver correctly decodes and sends back an ACK, or a maximum number of ARQ rounds (per message) is reached.

The authors are with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455, USA (email: pengwu@umn.edu; nihar@umn.edu).

ARQ provides an advantage by allowing for early termination once sufficient information has been received. As a result, it is most useful when there is considerably uncertainty in the amount/quality of information received. At the network layer, this might correspond to a setting where the network congestion is unknown to the transmitter. At the physical layer, which is the focus of this paper, this corresponds to a fading channel whose instantaneous quality is unknown to the transmitter.

Although H-ARQ is widely used in contemporary wireless systems such as HSPA [2], WiMax [3] (IEEE 802.16e) and 3GPP LTE [4], the majority of research on this topic has focused on code design, e.g., [5], [6], [7], while relatively little research has focused on performance analysis of H-ARQ [8]. Most relevant to the present work, in [9] Caire and Tuninetti established a relationship between H-ARQ throughput and mutual information in the limit of infinite block length. For multiple antenna systems, the diversity-multiplexing-delay tradeoff of H-ARQ was studied by El Gamal et al. [10], and the coding scheme achieving the optimal tradeoff was introduced; Chuang et al. [11] considered the optimal SNR exponent in the block-fading MIMO (multiple-input multiple-output) H-ARQ channel with discrete input signal constellation satisfying a short-term power constraint. H-ARQ has also been recently studied in quasi-static channels (i.e., the channel is fixed over all H-ARQ rounds) [12], [13] and shown to bring benefits to secrecy [14].

In this paper we build upon the results of [9] and perform a mutual information-based analysis of H-ARQ in block-fading channels. We consider a scenario where the fading is too fast to allow instantaneous channel quality feedback to the transmitter, and thus the transmitter only has knowledge of the channel statistics, but nonetheless each transmission experiences only a limited degree of channel selectivity. In this setting, rate adaptation can only be performed based on channel statistics and achieving a reasonable error/outage probability generally requires a conservative choice of rate if H-ARQ is not used. On the other hand, H-ARQ allows for *implicit* rate adaptation to the instantaneous channel quality because the receiver terminates transmission once the channel conditions experienced by a codeword are good enough to allow for decoding.

We analyze the long-term average transmitted rate achieved with H-ARQ, assuming that there is a maximum number of H-ARQ rounds and that a target outage probability at H-ARQ termination cannot be exceeded. We compare this rate to that achieved without H-ARQ in the same setting as well as to the ergodic capacity, which is the achievable rate in the idealized setting where instantaneous channel information is available to the transmitter. The main findings of the paper are that (a) H-ARQ generally provides a significant advantage over systems that do not use H-ARQ but have an equivalent level of channel selectivity, (b) the H-ARQ rate is reasonably close to the ergodic capacity in many practical

settings, and (c) the rate with H-ARQ is much less sensitive to the desired outage probability than an equivalent system that does not use H-ARQ.

The present work differs from prior literature in a number of important aspects. One key distinction is that we consider systems in which the rate is adapted to the average SNR such that a *constant* target outage probability is maintained at all SNR's, whereas most prior work has considered either fixed rate (and thus decreasing outage) [11] or increasing rate and *decreasing* outage as in the diversity-multiplexing tradeoff framework [10][15]. The fixed outage paradigm is consistent with contemporary wireless systems where an outage level near 1% is typical (see [16] for discussion), and certain conclusions depend heavily on the outage assumption. With respect to [9], note that the focus of [9] is on multi-user issues, e.g., whether or not a system becomes interference-limited at high SNR in the regime of very large delay, whereas we consider single-user systems and generally focus on performance with short delay constraints (i.e., maximum number of H-ARQ rounds). In addition, we use the *attempted* transmission rate, rather than the successful rate (which is used in [9]), as our performance metric. This is motivated by applications such as Voice-over-IP (VoIP), where a packet is dropped (and never retransmitted) if it cannot be decoded after the maximum number of H-ARQ rounds and quality of service is maintained by achieving the target (post-H-ARQ) outage probability. On the other hand, reliable data communication requires the use of higher-layer retransmissions whenever H-ARQ outages occur; in such a setting, the relevant metric is the successful rate, which is the product of the attempted transmission rate and the success probability (i.e., one minus the post-H-ARQ outage probability). If the post-H-ARQ outage is fixed to some target value (e.g., 1%), then studying the attempted rate is effectively equivalent to studying the successful rate.<sup>1</sup>

## II. SYSTEM MODEL

We consider a block-fading channel where the channel remains constant over a block but varies independently from one block to another. The  $t$ -th received symbol in the  $i$ -th block is given by:

$$y_{t,i} = \sqrt{\text{SNR}} h_i x_{t,i} + z_{t,i}, \quad (1)$$

where the index  $i = 1, 2, \dots$  indicates the block number,  $t = 1, 2, \dots, T$  indexes channel uses within a block, SNR is the average received SNR,  $h_i$  is the fading channel coefficient in the  $i$ -th block, and  $x_{t,i}$ ,  $y_{t,i}$ , and  $z_{t,i}$  are the transmitted symbol, received symbol, and additive noise, respectively. It is assumed that

<sup>1</sup>If the post-H-ARQ outage probability can be optimized, then a careful balancing between the attempted rate and higher-layer retransmissions should be conducted in order to maximize the successful rate. Although this is beyond the scope of the present paper, note that some results in this direction can be found in [17] [18].

$h_k$  is complex Gaussian (circularly symmetric) with unit variance and zero mean, and that  $h_1, h_2, \dots$  are i.i.d.. The noise  $z_{t,i}$  has the same distribution as  $h_k$  and is independent across channel uses and blocks. The transmitted symbol  $x_{t,i}$  is constrained to have unit average power; we consider Gaussian inputs, and thus  $x_{t,i}$  has the same distribution as the fading and the noise. Although we focus only on Rayleigh fading and single antenna systems, our basic insights can be extended to incorporate other fading distributions and MIMO as discussed in Section IV (Remark 1).

We consider the setting where the receiver has perfect channel state information (CSI), while the transmitter is aware of the channel distribution but does not know the instantaneous channel quality. This models a system in which the fading is too fast to allow for feedback of the instantaneous channel conditions from the receiver back to the transmitter, i.e., the channel coherence time is not much larger than the delay in the feedback loop. In cellular systems this is the case for moderate-to-high velocity users. This setting is often referred to as *open-loop* because of the lack of instantaneous channel tracking at the transmitter, although other forms of feedback, such as H-ARQ, are permitted. The relevant performance metrics, notably what we refer to as *outage probability* and fixed outage transmitted rate, are specified at the beginning of the relevant sections.

If H-ARQ is not used, we assume each codeword spans  $L$  fading blocks;  $L$  is therefore the channel selectivity experienced by each codeword. When H-ARQ is used, we make the following assumptions:

- The channel is constant within each H-ARQ round ( $T$  symbols), but is independent across H-ARQ rounds.<sup>2</sup>
- A maximum of  $M$  H-ARQ rounds are allowed. An outage is declared if decoding is not possible after  $M$  rounds, and this outage probability can be no larger than the constraint  $\epsilon$ .

Because the channel is assumed to be independent across H-ARQ rounds,  $M$  is the *maximum* amount of channel selectivity experienced by a codeword. When comparing H-ARQ and no H-ARQ, we set  $L = M$  such that maximum selectivity is equalized.

It is worth noting that these assumptions on the channel variation are quite reasonable for the fast-fading/open-loop scenarios. Transmission slots in modern systems are typically around one millisecond,

<sup>2</sup>An intuitive but somewhat misleading extension of the quasi-static fading model to the H-ARQ setting is to assume that the channel is constant for the duration of the H-ARQ rounds corresponding to a particular message/codeword, but is drawn independently across different messages. Because more H-ARQ rounds are needed to decode when the channel quality is poor, such a model actually changes the underlying fading distribution by increasing the probability of poor states and reducing the probability of good channel states. In this light, it is more accurate model the channel across H-ARQ rounds according to a stationary and ergodic random process with a high degree of correlation.

during which the channel is roughly constant even for fast fading.<sup>3</sup> An H-ARQ round generally corresponds to a single transmission slot, but subsequent ARQ rounds are separated in time by at least a few slots to allow for decoding and ACK/NACK feedback; thus the assumption of independent channels across H-ARQ rounds is reasonable. Moreover, a constraint on the number of H-ARQ rounds limits complexity (the decoder must retain information received in prior H-ARQ rounds in memory) and delay.

Throughout the paper we use the notation  $F_\epsilon^{-1}(X)$  to denote the solution  $y$  to the equation  $\mathbb{P}[X \leq y] = \epsilon$ , where  $X$  is a random variable; this quantity is well defined wherever it is used.

### III. PERFORMANCE WITHOUT H-ARQ: FIXED-LENGTH CODING

We begin by studying the baseline scenario where H-ARQ is not used and every codeword spans  $L$  fading blocks. In this setting the outage probability is the probability that mutual information received over the  $L$  fading blocks is smaller than the transmitted rate  $R$  [19, eq (5.83)]:

$$P_{\text{out}}(R, \text{SNR}) = \mathbb{P} \left[ \frac{1}{L} \sum_{i=1}^L \log_2(1 + \text{SNR}|h_i|^2) \leq R \right]. \quad (2)$$

where  $h_i$  is the channel in the  $i$ -th fading block. The outage probability reasonably approximates the decoding error probability for a system with strong coding [20] [21], and the achievability of this error probability has been rigorously shown in the limit of infinite block length ( $T \rightarrow \infty$ ) [22] [23].

Because the outage probability is a non-decreasing function of  $R$ , by setting the outage probability to  $\epsilon$  and solving for  $R$  we get the following straightforward definition of  $\epsilon$ -outage capacity [24]:

*Definition 1:* The  $\epsilon$ -outage capacity with outage constraint  $\epsilon$  and diversity order  $L$ , denoted by  $C_\epsilon^L(\text{SNR})$ , is the largest rate such that the outage probability in (2) is no larger than  $\epsilon$ :

$$C_\epsilon^L(\text{SNR}) \triangleq \max_{P_{\text{out}}(R, \text{SNR}) \leq \epsilon} R \quad (3)$$

Using notation introduced earlier, the  $\epsilon$ -outage capacity can be rewritten as

$$C_\epsilon^L(\text{SNR}) = F_\epsilon^{-1} \left( \frac{1}{L} \sum_{i=1}^L \log_2(1 + \text{SNR}|h_i|^2) \right) \quad (4)$$

$$= \log_2 \text{SNR} + F_\epsilon^{-1} \left( \frac{1}{L} \sum_{i=1}^L \log_2 \left( \frac{1}{\text{SNR}} + |h_i|^2 \right) \right). \quad (5)$$

For  $L = 1$ ,  $P_{\text{out}}(R, \text{SNR})$  can be written in closed form and inverted to yield  $C_\epsilon^1(\text{SNR}) = \log_2 \left( 1 + \log_e \left( \frac{1}{1-\epsilon} \right) \text{SNR} \right)$  [19]. For  $L > 1$  the outage probability cannot be written in closed form nor inverted, and therefore

<sup>3</sup>Frequency-domain channel variation within each H-ARQ round is briefly discussed in Section IV-B.

$C_\epsilon^L(\text{SNR})$  must be numerically computed. There are, however, two useful approximations to  $\epsilon$ -outage capacity. The first one is the high-SNR affine approximation [25], which adds a constant rate offset term to the standard multiplexing gain characterization.

*Theorem 1:* The high-SNR affine approximation to  $\epsilon$ -outage capacity is given by

$$C_\epsilon^L(\text{SNR}) = \log_2 \text{SNR} + F_\epsilon^{-1} \left( \frac{1}{L} \sum_{i=1}^L \log_2 (|h_i|^2) \right) + o(1), \quad (6)$$

where the notation implies that the  $o(1)$  term vanishes as  $\text{SNR} \rightarrow \infty$ .

*Proof:* The proof is identical to that of the high-SNR offset characterization of MIMO channels in [26, Theorem 1], noting that single antenna block fading is equivalent to a MIMO channel with a diagonal channel matrix. ■

In terms of standard high SNR notation where  $C(\text{SNR}) = \mathcal{S}_\infty(\log_2 \text{SNR} - \mathcal{L}_\infty) + o(1)$  [25][27], the multiplexing gain  $\mathcal{S}_\infty = 1$  and the rate offset  $\mathcal{L}_\infty = -F_\epsilon^{-1} \left( \frac{1}{L} \sum_{i=1}^L \log_2 (|h_i|^2) \right)$ . The rate offset is the difference between the  $\epsilon$ -outage capacity and the capacity of an AWGN channel with signal-to-noise ratio SNR. Although a closed form expression for  $\mathcal{L}_\infty$  cannot be found for  $L > 1$ , from [28],

$$\mathbb{P} \left[ \frac{1}{L} \sum_{i=1}^L \log_2 (|h_i|^2) \leq y \right] = 2^{yL} G_{1,L+1}^{L,1} (2^{yL} |_{0,0,\dots,0,-1}^0), \quad (7)$$

where  $G_{p,q}^{m,n} (x |_{b_1,\dots,b_q}^{a_1,\dots,a_p})$  is the Meijer G-function [29, eq. (9.301)]. Based on (6),  $\mathcal{L}_\infty$  therefore is the solution to  $2^{-\mathcal{L}_\infty L} G_{1,L+1}^{L,1} (2^{-\mathcal{L}_\infty L} |_{0,0,\dots,0,-1}^0) = \epsilon$ . The rate offset  $\mathcal{L}_\infty$  is plotted versus  $L$  in Fig. 1 for  $\epsilon = 0.01$ . As  $L \rightarrow \infty$  the offset converges to  $-\mathbb{E}[\log_2 (|h|^2)] \approx 0.83$ , the offset of the ergodic Rayleigh channel [27].

While the affine approximation is accurate at high SNR's, motivated by the Central Limit Theorem (CLT), an approximation that is more accurate for moderate and low SNR's is reached by approximating random variable  $\frac{1}{L} \sum_{i=1}^L \log_2(1 + \text{SNR}|h_i|^2)$  by a Gaussian random variable with the same mean and variance [30][31]. The mean  $\mu$  and variance  $\sigma^2$  of  $\log_2(1 + \text{SNR}|h|^2)$  are given by [32][33]:

$$\mu(\text{SNR}) = \mathbb{E}[\log_2(1 + \text{SNR}|h|^2)] = \log_2(e) e^{1/\text{SNR}} E_1(1/\text{SNR}), \quad (8)$$

$$\sigma^2(\text{SNR}) = \frac{2}{\text{SNR}} \log_2^2(e) e^{1/\text{SNR}} G_{3,4}^{4,0} \left( 1/\text{SNR} |_{0,-1,-1,-1}^{0,0,0} \right) - \mu^2(\text{SNR}), \quad (9)$$

where  $E_1(x) = \int_1^\infty t^{-1} e^{-xt} dt$ , and at high SNR the standard deviation  $\sigma(\text{SNR})$  converges to  $\frac{\pi \log_2 e}{\sqrt{6}}$  [33]. The mutual information is thus approximated by a  $\mathcal{N}(\mu(\text{SNR}), \frac{\sigma^2(\text{SNR})}{L})$ , and therefore

$$P_{\text{out}}(R, \text{SNR}) \approx Q \left( \frac{\sqrt{L}}{\sigma(\text{SNR})} (\mu(\text{SNR}) - R) \right), \quad (10)$$

where  $Q(\cdot)$  is the tail probability of a unit variance normal. Setting this quantity to  $\epsilon$  and then solving for  $R$  yields an  $\epsilon$ -outage capacity approximation [30, eq. (26)]:

$$C_\epsilon^L(\text{SNR}) \approx \mu(\text{SNR}) - \frac{\sigma(\text{SNR})}{\sqrt{L}} Q^{-1}(\epsilon). \quad (11)$$

The accuracy of this approximation depends on how accurately the CDF of a Gaussian matches the CDF (i.e., outage probability) of random variable  $\frac{1}{L} \sum_{i=1}^L \log_2(1 + \text{SNR}|h_i|^2)$ . In Fig. 2 both CDF's are plotted for  $L = 2$  and  $L = 10$ , and  $\text{SNR} = 0, 10$ , and  $20$  dB. As expected by the CLT, as  $L$  increases the approximation becomes more accurate. Furthermore, the match is less accurate for very small values of  $\epsilon$  because the tails of the Gaussian and the actual random variable do not precisely match. Finally, note that the match is not as accurate at low SNR's: this is because the mutual information random variable has a density close to a chi-square in this regime, and is thus not well approximated by a Gaussian. Although not accurate in all regimes, numerical results confirm that the Gaussian approximation is reasonably accurate for the range of interest for parameters (e.g.,  $0.01 \leq \epsilon \leq 0.2$  and  $0 \leq \text{SNR} \leq 20$  dB). More importantly, this approximation yields important insights.

In Fig. 3 the true  $\epsilon$ -outage capacity  $C_\epsilon^L(\text{SNR})$  and the affine and Gaussian approximations are plotted versus SNR for  $\epsilon = 0.01$  and  $L = 3, 10$ . The Gaussian approximation is reasonably accurate at moderate SNR's, and is more accurate for larger values of  $L$ . On the other hand, the affine approximation, which provides a correct high SNR offset, is asymptotically tight at high SNR.

### A. Ergodic Capacity Gap

When evaluating the effect of the diversity order  $L$ , it is useful to compare the ergodic capacity  $\mu(\text{SNR})$  and  $C_\epsilon^L(\text{SNR})$ . By Chebyshev's inequality, for any  $0 < w < \mu$ ,

$$\mathbb{P} \left[ \left| \mu(\text{SNR}) - \frac{1}{L} \sum_{i=1}^L \log_2(1 + \text{SNR}|h_i|^2) \right| \geq w \right] \leq \frac{\sigma^2(\text{SNR})}{Lw^2} \quad (12)$$

By replacing  $w$  with  $\mu(\text{SNR}) - R$  and equating the right hand side (RHS) with  $\epsilon$ , we get

$$\mu(\text{SNR}) - \frac{\sigma(\text{SNR})}{\sqrt{L\epsilon}} \leq C_\epsilon^L(\text{SNR}) \leq \mu(\text{SNR}) + \frac{\sigma(\text{SNR})}{\sqrt{L\epsilon}}. \quad (13)$$

This implies  $C_\epsilon^L(\text{SNR}) \rightarrow \mu(\text{SNR})$  as  $L \rightarrow \infty$ , as intuitively expected; reasonable values of  $\epsilon$  are smaller than 0.5, and thus we expect convergence to occur from below.

In order to capture the speed at which this convergence occurs, we define the quantity  $\Delta_{\text{EC-FD}}$  as the difference between the ergodic and  $\epsilon$ -outage capacities. Based on (13) we can upper bound  $\Delta_{\text{EC-FD}}$  as:

$$\Delta_{\text{EC-FD}}(\text{SNR}) = \mu(\text{SNR}) - C_\epsilon^L(\text{SNR}) \leq \frac{\sigma(\text{SNR})}{\sqrt{L\epsilon}}. \quad (14)$$

This bound shows that the rate gap goes to zero at least as fast as  $O(1/\sqrt{L})$ . Although we cannot rigorously claim that  $\Delta_{\text{EC-FD}}$  is of order  $1/\sqrt{L}$ , by (11) the Gaussian approximation to this quantity is:

$$\Delta_{\text{EC-FD}}(\text{SNR}) \approx \frac{\sigma(\text{SNR})}{\sqrt{L}} Q^{-1}(\epsilon), \quad (15)$$

which is also  $O(1/\sqrt{L})$ . This approximation becomes more accurate as  $L \rightarrow \infty$ , by the CLT, and thus is expected to correctly capture the scaling with  $L$ . Note that (15) has the interpretation that the rate must be  $\frac{Q^{-1}(\epsilon)}{\sqrt{L}}$  deviations below the ergodic capacity  $\mu(\text{SNR})$  in order to ensure  $1 - \epsilon$  reliability. In Fig. 4 the actual capacity gap and the approximation in (15) are plotted for  $\epsilon = 0.01$  and  $\epsilon = 0.05$  with  $\text{SNR} = 20$  dB, and a reasonable match between the approximation and the exact gap is seen.

#### IV. PERFORMANCE WITH HYBRID-ARQ

We now move on to the analysis of hybrid-ARQ, which will be shown to provide a significant performance advantage relative to the baseline of non-H-ARQ performance. H-ARQ is clearly a variable-length code, in which case the average transmission rate must be suitably defined. If each message contains  $b$  information bits and each ARQ round corresponds to  $T$  channel symbols, then the *initial transmission rate* is  $R_{\text{init}} \triangleq \frac{b}{T}$  bits/symbol. If random variable  $X_i$  denotes the number of H-ARQ rounds used for the  $i$ -th message, then a total of  $\sum_{i=1}^N X_i$  H-ARQ rounds are used and the average transmission rate (in bits/symbol or bps/Hz) across those  $N$  messages is:

$$\frac{Nb}{T \sum_{i=1}^N X_i} = \frac{R_{\text{init}}}{\frac{1}{N} \sum_{i=1}^N X_i}. \quad (16)$$

We are interested in the long-term average transmission rate, i.e., the case where  $N \rightarrow \infty$ . By the law of large numbers (note that the  $X_i$ 's are i.i.d. in our model),  $\frac{1}{N} \sum_{i=1}^N X_i \rightarrow \mathbb{E}[X]$  and thus the rate converges to

$$\frac{R_{\text{init}}}{\mathbb{E}[X]} \text{ bits/symbol} \quad (17)$$

Here  $X$  is the random variable representing the number of H-ARQ rounds per message; this random variable is determined by the specifics of the H-ARQ protocol.

In the remainder of the paper we focus on incremental redundancy (IR) H-ARQ because it is the most powerful type of H-ARQ, although we compare IR to Chase combining in Section IV-E. In [9] it is shown that mutual information is *accumulated* over H-ARQ rounds when IR is used, and that decoding is possible once the accumulated mutual information is larger than the number of information bits in the message. Therefore, the number of H-ARQ rounds  $X$  is the smallest number  $m$  such that:

$$\sum_{i=1}^m \log_2(1 + \text{SNR}|h_i|^2) > R_{\text{init}}. \quad (18)$$

The number of rounds is upper bounded by  $M$ , and an outage occurs whenever the mutual information after  $M$  rounds is smaller than  $R_{\text{init}}$ :

$$P_{\text{out}}^{\text{IR},M}(R_{\text{init}}) = \mathbb{P} \left[ \sum_{i=1}^M \log_2(1 + \text{SNR}|h_i|^2) \leq R_{\text{init}} \right]. \quad (19)$$

This is the same as the expression for outage probability of  $M$ -order diversity without H-ARQ in (2), except that mutual information is summed rather than averaged over the  $M$  rounds. This difference is a consequence of the fact that  $R_{\text{init}}$  is defined for transmission over one round rather than all  $M$  rounds; dividing by  $\mathbb{E}[X]$  in (17) to obtain the average transmitted rate makes the expressions consistent. Due to this relationship, if the initial rate is set as  $R_{\text{init}} = M \cdot C_\epsilon^M$ , where  $C_\epsilon^M$  is the  $\epsilon$ -outage capacity for  $M$ -order diversity without H-ARQ, then the outage at H-ARQ termination is  $\epsilon$ .

In order to simplify expressions, it is useful to define  $A_k(R_{\text{init}})$  as the probability that the accumulated mutual information after  $k$  rounds is smaller than  $R_{\text{init}}$ :

$$A_k(R_{\text{init}}) \triangleq \mathbb{P} \left[ \sum_{i=1}^k \log_2(1 + \text{SNR}|h_i|^2) \leq R_{\text{init}} \right]. \quad (20)$$

The expected number of H-ARQ rounds per message is therefore given by:

$$\mathbb{E}[X] = 1 + \sum_{k=1}^{M-1} \mathbb{P}[X > k] = 1 + \sum_{k=1}^{M-1} A_k(R_{\text{init}}). \quad (21)$$

The long-term average transmitted rate, which is denoted as  $C_\epsilon^{\text{IR},M}$ , is defined by (17). With initial rate  $R_{\text{init}} = M \cdot C_\epsilon^M$  we have:<sup>4</sup>

$$C_\epsilon^{\text{IR},M} \triangleq \frac{R_{\text{init}}}{\mathbb{E}[X]} = \left( \frac{M}{\mathbb{E}[X]} \right) C_\epsilon^M. \quad (22)$$

Note that  $C_\epsilon^{\text{IR},M}$  is the *attempted* long-term average transmission rate, as discussed in Section I. For the sake of brevity this quantity is referred to as the H-ARQ rate; this is not to be confused with the initial rate  $R_{\text{init}}$ . Similarly, we refer to  $\epsilon$ -outage capacity  $C_\epsilon^M$  as the non-H-ARQ rate in the rest of the paper.

Because  $\mathbb{E}[X] \leq M$ , the H-ARQ rate is at least as large as the non-H-ARQ rate, i.e.,  $C_\epsilon^{\text{IR},M} \geq C_\epsilon^M$ , and the advantage with respect to the non-H-ARQ benchmark is precisely the multiplicative factor  $\frac{M}{\mathbb{E}[X]}$ . This difference is explained as follows. Because  $R_{\text{init}} = M \cdot C_\epsilon^M$ , each message/packet contains  $C_\epsilon^M MT$  information bits regardless of whether H-ARQ is used. Without H-ARQ these bits are always transmitted over  $MT$  symbols, whereas with H-ARQ an average of only  $\mathbb{E}[X]T$  symbols are required.

<sup>4</sup>All quantities in this expression except  $M$  are actually functions of SNR. For the sake of compactness, however, dependence upon SNR is suppressed in this and subsequent expressions, except where explicit notation is necessary.

In Fig. 5 the average rates with ( $C_\epsilon^{\text{IR},M}$ ) and without H-ARQ ( $C_\epsilon^M$ ) are plotted versus SNR for  $\epsilon = 0.01$  and  $M = 1, 2$  and  $6$  ( $M = 1$  does not allow for H-ARQ in our model). Ergodic capacity is also plotted as a reference. Based on the figure, we immediately notice:

- H-ARQ with 6 rounds outperforms H-ARQ with 2 rounds.
- H-ARQ provides a significant advantage relative to non-H-ARQ for the same value of  $M$  for a wide range of SNR's, but this advantage vanishes at high SNR.

Increasing rate with  $M$  is to be expected, because larger  $M$  corresponds to more time diversity and more early termination opportunities. The behavior with respect to SNR is perhaps less intuitive. The remainder of this section is devoted to quantifying and explaining the behavior seen in Fig. 5. We begin by extending the Gaussian approximation to H-ARQ, then examine performance scaling with respect to  $M$ , SNR, and  $\epsilon$ , and finally compare IR to Chase Combining.

#### A. Gaussian Approximation

By the definition of  $A_k(\cdot)$  and (21)-(22), the H-ARQ rate can be written as:

$$C_\epsilon^{\text{IR},M} = \frac{A_M^{-1}(\epsilon)}{1 + \sum_{k=1}^{M-1} A_k(A_M^{-1}(\epsilon))}, \quad (23)$$

where  $A_M^{-1}(\cdot)$  refers to the inverse of function  $A_M(\cdot)$ . If we use the approach of Section III and approximate the mutual information accumulated in  $k$  rounds by a Gaussian with mean  $\mu k$  and variance  $\sigma^2 k$ , where  $\mu$  and  $\sigma^2$  are defined in (8) and (9), we have:

$$A_k(R_{\text{init}}) \approx Q\left(\frac{\mu k - R_{\text{init}}}{\sigma\sqrt{k}}\right). \quad (24)$$

Similar to (11), the initial rate  $R_{\text{init}} = A_M^{-1}(\epsilon)$  can be approximated as  $M\left[\mu - \frac{\sigma}{\sqrt{M}}Q^{-1}(\epsilon)\right]$ . Applying the approximation of  $A_k(R_{\text{init}})$  to each term in (23) and using the property  $1 - Q(x) = Q(-x)$  yields:

$$C_\epsilon^{\text{IR},M} \approx \frac{M\left[\mu - \frac{\sigma}{\sqrt{M}}Q^{-1}(\epsilon)\right]}{M - \sum_{k=1}^{M-1} Q\left(\frac{M-k}{\sqrt{k}}\frac{\mu}{\sigma} - \sqrt{\frac{M}{k}}Q^{-1}(\epsilon)\right)}. \quad (25)$$

This approximation is easier to compute than the actual H-ARQ rate and is reasonably accurate. Furthermore, it is useful for the insights it can provide.

#### B. Scaling with H-ARQ Rounds $M$

In this section we study the dependence of the H-ARQ rate on  $M$ . We first show convergence to the ergodic capacity as  $M \rightarrow \infty$ :

*Theorem 2:* For any SNR, the H-ARQ rate converges to the ergodic capacity as  $M \rightarrow \infty$ :

$$\lim_{M \rightarrow \infty} C_\epsilon^{\text{IR},M}(\text{SNR}) = \mu(\text{SNR}) \quad (26)$$

*Proof:* See Appendix A. ■

To quantify how fast this convergence is, similar to Section III-A we investigate the difference between the ergodic capacity and the H-ARQ rate. Defining  $\Delta_{\text{EC-IR}} \triangleq \mu(\text{SNR}) - C_\epsilon^{\text{IR},M}(\text{SNR})$  we have

$$\Delta_{\text{EC-IR}} = \frac{\mu}{\mathbb{E}[X]} \left( \mathbb{E}[X] - \frac{M \cdot C_\epsilon^M}{\mu} \right) \approx \frac{\mu}{\mathbb{E}[X]} \left( \mathbb{E}[X] - \left( M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} \right) \right) \quad (27)$$

where the approximation follows from  $C_\epsilon^M \approx \mu - \frac{\sigma Q^{-1}(\epsilon)}{\sqrt{M}}$  in (11). Because  $\mathbb{E}[X]$  is on the order of  $M$  (as established in the proof of Theorem 2), the key is the behavior of the term  $\mathbb{E}[X] - \left( M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} \right)$ .

To better understand  $\mathbb{E}[X]$  we again return to the Gaussian approximation. While the CDF of  $X$  is defined by  $\mathbb{P}(X \leq k) = 1 - A_k(R_{\text{init}})$  (for  $k = 1, \dots, M-1$ ), we use  $\tilde{X}$  to denote the random variable using the Gaussian approximation and thus define its CDF (for integers  $k$ ) as:

$$\mathbb{P}(\tilde{X} \leq k) = Q \left( \frac{\mu(M-k) - \sqrt{M}\sigma Q^{-1}(\epsilon)}{\sigma\sqrt{k}} \right) \quad (28)$$

where we have used  $A_k(R_{\text{init}}) \approx Q \left( \frac{\mu k - R_{\text{init}}}{\sigma\sqrt{k}} \right)$  evaluated with  $R_{\text{init}} = MC_\epsilon^M \approx M\mu - \sqrt{M}\sigma Q^{-1}(\epsilon)$ . From this expression, we can immediately see that the *median* of  $\tilde{X}$  is  $\left\lceil M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} \right\rceil$  [34]. If this was equal to the mean of  $X$ , then by (27) the rate difference would be well approximated by  $\frac{\beta\mu}{M}$ , where  $\beta$  is the difference between  $\left\lceil M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} \right\rceil$  and  $\left( M - \frac{\beta}{\mu} Q^{-1}(\epsilon) \sqrt{M} \right)$  and thus is no larger than one. By studying the characteristics of  $\tilde{X}$  (and of  $X$ ) we can see that the median is in fact quite close to the mean. A tedious calculation in Appendix B gives the following approximation to  $\mathbb{E}[X]$ :

$$\mathbb{E}[X] \approx M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} + 0.5(1-\epsilon) - \frac{\sigma}{\mu} \sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx, \quad (29)$$

which is reasonably accurate for large  $M$ . The most important factor is the term  $0.5(1-\epsilon)$ , which is due to the fact that only an integer number of H-ARQ rounds can be used. The factor  $-\frac{\sigma}{\mu} \sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx$  exists because the random variable is truncated at the point where its CDF is  $1-\epsilon$ .

Applying this into (27), the rate difference can be approximated as:

$$\Delta_{\text{EC-IR}} \approx \frac{\mu \left( 0.5(1-\epsilon) - \frac{\sigma}{\mu} \sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx \right)}{M - \frac{\sigma}{\mu} Q^{-1}(\epsilon) \sqrt{M} - \frac{\sigma}{\mu} \sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx + 0.5(1-\epsilon)}. \quad (30)$$

The denominator increases with  $M$  at the order of  $M$  (more precisely as  $M - \sqrt{M}$ ), while the numerator actually decreases with  $M$  and can even become negative if  $M$  is extremely large. For reasonable values of  $M$ , however, the negative term in the numerator is essentially inconsequential (for example, if  $\epsilon = 0.01$

and  $\text{SNR} = 10$  dB, the negative term is much smaller than  $0.5(1 - \epsilon)$  for  $M < 5000$ ) and thus can be reasonably neglected. By ignoring this negative term and replacing the denominator with the leading order  $M$  term, we get a further approximation of the rate gap:

$$\Delta_{\text{EC-IR}} \approx \frac{0.5(1 - \epsilon)\mu}{M} \quad (31)$$

Based on this approximation, we see that the rate gap decreases roughly on the order  $O(1/M)$ , rather than the  $O(1/\sqrt{M})$  decrease without H-ARQ. In Fig. 6 we plot the exact capacity gap with and without H-ARQ, as well as the Gaussian approximation to the H-ARQ gap (30) and its simplified form in (31) for  $\epsilon = 0.01$  at  $\text{SNR} = 10$  dB. Both approximations are seen to be reasonably accurate especially for large  $M$ . In the inset plot, which is in log-log scale, we see that the exact capacity gap goes to zero at order  $1/M$ , consistent with the result obtained from our approximation.

The fast convergence with H-ARQ can be intuitively explained as follows. If transmission could be stopped precisely when enough mutual information has been received, the transmitted rate would be exactly matched to the instantaneous mutual information and thus ergodic capacity would be achieved. When H-ARQ is used, however, transmission can only be terminated at the end of a round as, opposed to within a round, and thus a small amount of the transmission can be wasted. This "rounding error", which is reflected in the  $0.5(1 - \epsilon)$  term in (30) and (31), is essentially the only penalty incurred by using H-ARQ rather than explicit rate adaptation.

**Remark 1:** Because the value of H-ARQ depends primarily on the mean and variance of the mutual information in each H-ARQ round, our basic insights can be extended to multiple-antenna channels and to channels with frequency (or time) diversity within each ARQ round if the change in mean and variance is accounted for. For example, with order  $F$  frequency diversity the mutual information in the  $i$ -th H-ARQ round becomes  $\frac{1}{F} \sum_{l=1}^F \log(1 + \text{SNR}|h_{i,l}|^2)$ , where  $h_{i,l}$  is the channel in the  $i$ -th round on the  $l$ -th frequency channel. The mean mutual information is unaffected, while the variance is decreased by a factor of  $\frac{1}{F}$ .  $\diamond$

### C. Scaling with SNR

In this section we quantify the behavior of H-ARQ as a function of the average SNR.<sup>5</sup> Fig. 5 indicated that the benefit of H-ARQ vanishes at high SNR, and the following theorem makes this precise:

<sup>5</sup>Because constant outage corresponds to the full-multiplexing point, the results of [10] imply that  $C_e^{\text{IR},M}$  cannot have a multiplexing gain/pre-log larger than one (in [10] it is shown that H-ARQ does not increase the full multiplexing point). However, the DMT-based results of [10] do not provide rate-offset characterization as in Theorems 3 and 4.

*Theorem 3:* If SNR is taken to infinity while keeping  $M$  fixed, the expected number of H-ARQ rounds converges to  $M$  and the H-ARQ rate converges to  $C_\epsilon^M(\text{SNR})$ , the non-H-ARQ rate with the same selectivity:

$$\lim_{\text{SNR} \rightarrow \infty} \mathbb{E}[X] = M \quad (32)$$

$$\lim_{\text{SNR} \rightarrow \infty} [C_\epsilon^{\text{IR},M}(\text{SNR}) - C_\epsilon^M(\text{SNR})] = 0 \quad (33)$$

*Proof:* See Appendix C. ■

The intuition behind this result can be gathered from Fig. 7, where the CDF's of the accumulated mutual information after 2 and 3 rounds are plotted for for SNR = 10 and 40 dB. If  $M = 3$  the initial rate is set at the  $\epsilon$ -point of the CDF of  $\sum_{i=1}^3 \log(1 + \text{SNR}|h_i|^2)$ . Because the CDF's overlap for  $M = 2$  and  $M = 3$  considerably when SNR = 10 dB, there is a large probability that sufficient mutual information is accumulated after 2 rounds and thus early termination occurs. However, the overlap between these CDF's disappears as SNR increases, because  $\sum_{i=1}^k \log(1 + \text{SNR}|h_i|^2) \approx k \log \text{SNR} + \sum_{i=1}^k \log(|h_i|^2)$ , and thus the early termination probability vanishes.

Although the H-ARQ advantage eventually vanishes, the advantage persists throughout a large SNR range and the Gaussian approximation (Section IV-A) can be used to quantify this. The probability of terminating in strictly less than  $M$  rounds is approximated by:

$$\mathbb{P}[X \leq M - 1] \approx Q\left(\frac{\mu(\text{SNR}) - \sqrt{M}\sigma(\text{SNR})Q^{-1}(\epsilon)}{\sigma\sqrt{M-1}}\right) \quad (34)$$

In order for this approximation to be greater than one-half we require the numerator inside the  $Q$ -function to be less than zero, which corresponds to

$$\frac{\mu(\text{SNR})}{\sigma(\text{SNR})} \leq \sqrt{M}Q^{-1}(\epsilon) \quad \text{or} \quad \sqrt{M} \geq \frac{\mu(\text{SNR})}{\sigma(\text{SNR})Q^{-1}(\epsilon)} \quad (35)$$

As SNR increases  $\mu(\text{SNR})$  increases without bound whereas  $\sigma(\text{SNR})$  converges to a constant. Thus  $\mu(\text{SNR})/\sigma(\text{SNR})$  increases quickly with SNR, which makes the probability of early termination vanish. From this we see that the H-ARQ advantage lasts longer (in terms of SNR) when  $M$  is larger. The second inequality in (35) captures an alternative viewpoint, which is roughly the minimum value of  $M$  required for H-ARQ to provide a significant advantage.

Motivated by naive intuition that the H-ARQ rate is monotonically increasing in the initial rate  $R_{\text{init}}$ , up to this point we have chosen  $R_{\text{init}} = MC_\epsilon^M = A_M^{-1}(\epsilon)$  such that outage at H-ARQ termination is exactly  $\epsilon$ . However, it turns out that the H-ARQ rate is not always monotonic in  $R_{\text{init}}$ . In Fig. 8, the H-ARQ rate  $\frac{R_{\text{init}}}{\mathbb{E}[X]}$  is plotted versus initial rate  $R_{\text{init}}$  for  $M = 2, 3$ , and 4 for SNR = 10 dB (left) and 30 dB (right). At

10 dB,  $\frac{R_{\text{init}}}{\mathbb{E}[X]}$  monotonically increases with  $R_{\text{init}}$  and thus there is no advantage to optimizing the initial rate. At 30 dB, however,  $\frac{R_{\text{init}}}{\mathbb{E}[X]}$  behaves non-monotonically with  $R_{\text{init}}$ . We therefore define  $\tilde{C}_\epsilon^{\text{IR},M}(\text{SNR})$  as the maximized H-ARQ rate, where the maximization is performed over all values of initial rate  $R_{\text{init}}$  such that the outage constraint  $\epsilon$  is not violated:

$$\tilde{C}_\epsilon^{\text{IR},M}(\text{SNR}) \triangleq \max_{R_{\text{init}} \leq A_M^{-1}(\epsilon)} \frac{R_{\text{init}}}{\mathbb{E}[X]} \quad (36)$$

The local maxima seen in Fig. 8 appear to preclude a closed form solution to this maximization. Although optimization of the initial rate provides an advantage over a certain SNR range, the following theorem shows that it does not provide an improvement in the high-SNR offset:

*Theorem 4:* H-ARQ with an optimized initial rate, i.e.,  $\tilde{C}_\epsilon^{\text{IR},M}(\text{SNR})$ , achieves the same high-SNR offset as unoptimized H-ARQ  $C_\epsilon^{\text{IR},M}(\text{SNR})$

$$\lim_{\text{SNR} \rightarrow \infty} \left[ \tilde{C}_\epsilon^{\text{IR},M}(\text{SNR}) - C_\epsilon^{\text{IR},M}(\text{SNR}) \right] = 0 \quad (37)$$

Furthermore, the only initial rate (ignoring  $o(1)$  terms) that achieves the correct offset is the unoptimized value  $R_{\text{init}} = MC_\epsilon^M$ .

*Proof:* See Appendix D. ■

In Fig. 9, rates with and without optimization of the initial rate are plotted for  $\epsilon = 0.01$  and  $M = 2, 6$ . For  $M = 2$  optimization begins to make a difference at the point where the unoptimized curve abruptly decreases towards  $C_\epsilon^M$  around 25 dB, but this advantage vanishes around 55 dB. For  $M = 6$  the advantage of initial rate optimization comes about at a much higher SNR, consistent with (35). Convergence of  $\tilde{C}_\epsilon^{\text{IR},6}(\text{SNR})$  to  $C_\epsilon^{\text{IR},6}(\text{SNR})$  does eventually occur, but is not visible in the figure.

#### D. Scaling with Outage Constraint $\epsilon$

Another advantage of H-ARQ is that the H-ARQ rate is generally less sensitive to the desired outage probability  $\epsilon$  than an equivalent non-H-ARQ system. This advantage is clearly seen in Fig. 10, where the H-ARQ and non-H-ARQ rates are plotted versus  $\epsilon$  for  $M = 5$  at SNR = 0, 10 and 20 dB. When  $\epsilon$  is large (e.g., roughly around 0.5) H-ARQ provides almost no advantage: a large outage corresponds to a large initial rate, which in turn means early termination rarely occurs. However, for more reasonable values of  $\epsilon$ , the H-ARQ rate is roughly constant with respect to  $\epsilon$  whereas the non-H-ARQ rate decreases sharply as  $\epsilon \rightarrow 0$ . The transmitted rate must be decreased in order to achieve a smaller  $\epsilon$  (with or without H-ARQ), but with H-ARQ this decrease is partially compensated by the accompanying decreasing in the number of rounds  $\mathbb{E}[X]$ .

### E. Chase Combining

If Chase combining is used, a packet is retransmitted whenever a NACK is received and the receiver performs maximal-ratio-combining (MRC) on all received packets. As a result, SNR rather the mutual information is accumulated over H-ARQ rounds and the outage probability is given by:

$$P_{\text{out}}^{\text{CC},M}(R_{\text{init}}) = \mathbb{P} \left[ \log_2 \left( 1 + \text{SNR} \sum_{i=1}^M |h_i|^2 \right) \leq R_{\text{init}} \right]. \quad (38)$$

For outage  $\epsilon$ , the initial rate is  $R_{\text{init}} = \log_2 \left( 1 + F_{\epsilon}^{-1} \left( \sum_{i=1}^M |h_i|^2 \right) \text{SNR} \right)$ .

Different from IR, the expected number of H-ARQ rounds in CC is not dependent on SNR and thus the average rate for outage  $\epsilon$  can be written in closed form:

$$C_{\epsilon}^{\text{CC},M}(\text{SNR}) = \frac{R_{\text{init}}}{\mathbb{E}[X]} = \frac{\log_2 \left( 1 + F_{\epsilon}^{-1} \left( \sum_{i=1}^M |h_i|^2 \right) \text{SNR} \right)}{M - e^{-F_{\epsilon}^{-1} \left( \sum_{i=1}^M |h_i|^2 \right)} \sum_{k=1}^{M-1} (M-k) \frac{(F_{\epsilon}^{-1} \left( \sum_{i=1}^M |h_i|^2 \right))^{k-1}}{(k-1)!}}, \quad (39)$$

where the denominator is  $\mathbb{E}[X]$ . According to (39), we can get the high SNR affine approximation as:

$$C_{\epsilon}^{\text{CC},M}(\text{SNR}) = \frac{1}{\mathbb{E}[X]} \log_2 \text{SNR} + \frac{1}{\mathbb{E}[X]} \log_2 \left( F_{\epsilon}^{-1} \left( \sum_{i=1}^M |h_i|^2 \right) \right) + o(1). \quad (40)$$

Because  $\mathbb{E}[X] > 1$  for any positive outage value, the pre-log factor (i.e., multiplexing gain) is  $\frac{1}{\mathbb{E}[X]}$  and thus is less than one. This implies that CC performs poorly at high SNR. This is to be expected because CC is essentially a repetition code, which is spectrally inefficient at high SNR. As with IR, the performance of CC at high SNR can be improved through rate optimization. At high SNR, the pre-log is critical and thus the initial rate should be selected so that  $\mathbb{E}[X]$  is close to one and thereby avoiding H-ARQ altogether. Even with optimization, CC is far inferior to IR at moderate and high SNR's. On the other hand, CC performs reasonably well at low SNR. This is because  $\log(1+x) \approx x$  for small values of  $x$ , and thus SNR-accumulation is nearly equivalent to mutual information-accumulation. In Fig. 11 rate-optimized IR and CC are plotted for  $M = 2, 4$  and  $\epsilon = 0.01$ , and the results are consistent with the above intuitions.

## V. CONCLUSION

In this paper we have studied the performance of hybrid-ARQ in the context of an open-loop/fast-fading system in which the transmission rate is adjusted as a function of the average SNR such that a target outage probability is not exceeded. The general findings are that H-ARQ provides a significant rate advantage relative to a system not using H-ARQ at reasonable SNR levels, and that H-ARQ provides a rate quite close to the ergodic capacity even when the channel selectivity is limited.

There appear to be some potentially interesting extensions of this work. Contemporary cellular systems utilize simple ARQ on top of H-ARQ, and it is not fully understood how to balance these reliability mechanisms; some results in this direction are presented in [17]. Although we have assumed error-free ACK/NACK feedback, such errors can be quite important (c.f., [35]) and merit further consideration. Finally, while we have considered only the mutual information of Gaussian inputs, it is of interest to extend the results to discrete constellations and possibly compare to the performance of actual codes.

## VI. ACKNOWLEDGMENTS

The authors gratefully acknowledge Prof. Robert Heath for suggesting use of the Gaussian approximation, and Prof. Gerhard Wunder for suggesting use of Chebyshev's inequality in Section III-A.

### APPENDIX A

#### PROOF OF THEOREM 2

Because  $C_\epsilon^{\text{IR},M} = \frac{M}{\mathbb{E}[X]} C_\epsilon^M$  and  $\lim_{M \rightarrow \infty} C_\epsilon^M = \mu$  (Section III-A), we can prove  $\lim_{M \rightarrow \infty} C_\epsilon^{\text{IR},M} = \mu$  by showing  $\lim_{M \rightarrow \infty} \frac{M}{\mathbb{E}[X]} = 1$ . Because  $\mathbb{E}[X] \leq M$ , we can show this simply by showing  $\mathbb{E}[X]$  is of order  $M$ . For notational convenience we define  $Y_i \triangleq \log_2(1 + \text{SNR}|h_i|^2)$ , and then have:

$$\mathbb{E}[X] \geq 1 + \sum_{k=1}^{M-1} \mathbb{P} \left[ \sum_{i=1}^k Y_i \leq M \left( \mu - \frac{\sigma}{\sqrt{M\epsilon}} \right) \right] \quad (41)$$

$$\geq 1 + \sum_{k=1}^{\lceil M - M^{\frac{3}{4}} \rceil} \mathbb{P} \left[ \sum_{i=1}^k Y_i \leq M\mu - \sigma \sqrt{\frac{M}{\epsilon}} \right] \quad (42)$$

$$\geq 1 + \lceil M - M^{\frac{3}{4}} \rceil \cdot \mathbb{P} \left[ \sum_{i=1}^{\lceil M - M^{\frac{3}{4}} \rceil} Y_i \leq M\mu - \sigma \sqrt{\frac{M}{\epsilon}} \right], \quad (43)$$

where the first line holds because  $\mathbb{E}[X]$  is increasing in  $R_{\text{init}}$  and  $R_{\text{init}} = MC_\epsilon^M \geq M \left( \mu - \frac{\sigma}{\sqrt{M\epsilon}} \right)$  from (13), the second holds because the summands are non-negative, and the last line because the summands are decreasing in  $k$ . A direct application of the CLT shows  $\mathbb{P} \left[ \sum_{i=1}^{\lceil M - M^{\frac{3}{4}} \rceil} Y_i \leq M\mu - \sigma \sqrt{\frac{M}{\epsilon}} \right] \rightarrow 1$  as  $M \rightarrow \infty$ , and thus, with some straightforward algebra, we have  $\lim_{M \rightarrow \infty} \frac{M}{\mathbb{E}[X]} \rightarrow 1$ .

### APPENDIX B

#### PROOF OF (29)

Firstly, we relax the constraint on  $\tilde{X}$  (discreteness and finiteness) to define a new continuous random variable  $\hat{X}$ , which is distributed along the whole real line. The CDF of  $\hat{X}$  (for all real  $x$ ) is

$$\mathbb{P} \left( \hat{X} \leq x \right) = Q \left( \frac{\mu(M-x) - \sqrt{M}\sigma Q^{-1}(\epsilon)}{\sigma\sqrt{x}} \right) \quad (44)$$

Now if we consider the distribution of  $\hat{X} - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right)$ , we have

$$\mathbb{P}\left[\hat{X} - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right) \leq x\right] = Q\left(\frac{-\mu x}{\sigma\sqrt{M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M} + x}}\right) \quad (45)$$

where the equality follows from (44). Notice as  $M \rightarrow \infty$ ,  $\sqrt{M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M} + x} \rightarrow \sqrt{M}$ , so

$$Q\left(\frac{-\mu x}{\sigma\sqrt{M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M} + x}}\right) \rightarrow Q\left(\frac{-\mu x}{\sigma\sqrt{M}}\right) = \Phi\left(\frac{x}{\frac{\sigma}{\mu}\sqrt{M}}\right), \quad (46)$$

where  $\Phi(\cdot)$  is the standard normal CDF with zero mean and unit variance. So as  $M \rightarrow \infty$ , the limiting distribution of  $\hat{X}$ , denoted by  $\hat{\Phi}(\cdot)$ , goes to  $\mathcal{N}\left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}, \left(\frac{\sigma}{\mu}\right)^2 M\right)$ , which is

$$\hat{\Phi}(x) = \Phi\left(\frac{x - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right)}{\frac{\sigma}{\mu}\sqrt{M}}\right), \quad \text{for } x \in \mathbb{R} \quad (47)$$

Since  $\mathbb{E}[\tilde{X}]$  is an approximation to  $\mathbb{E}[X]$ , then we focus on evaluating  $\mathbb{E}[\tilde{X}]$  for large  $M$ :

$$\begin{aligned} \mathbb{E}[\tilde{X}] &= \sum_{k=0}^{M-1} \left(1 - \mathbb{P}[\tilde{X} \leq k]\right) = M - \sum_{k=1}^{M-1} \tilde{\Phi}(k) \stackrel{(a)}{=} M - \sum_{k=1}^{M-1} \hat{\Phi}(k) \\ &\approx M - \left(\int_1^M \hat{\Phi}(x) dx - \sum_{k=1}^{M-1} \frac{\hat{\Phi}(k+1) - \hat{\Phi}(k)}{2}\right) \\ &\stackrel{(b)}{\approx} M - \left(\int_1^M \Phi\left(\frac{x - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right)}{\frac{\sigma}{\mu}\sqrt{M}}\right) dx - 0.5(1 - \epsilon)\right) \\ &= M - \int_{M - 2\frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}}^M \Phi\left(\frac{x - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right)}{\frac{\sigma}{\mu}\sqrt{M}}\right) dx - \\ &\quad \int_1^{M - 2\frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}} \Phi\left(\frac{x - \left(M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}\right)}{\frac{\sigma}{\mu}\sqrt{M}}\right) dx + 0.5(1 - \epsilon) \end{aligned} \quad (48)$$

where (a) holds since  $\tilde{\Phi}(k) = \hat{\Phi}(k)$  when  $k = 1, 2, \dots, M - 1$  and (b) follows from  $\hat{\Phi}(M) = 1 - \epsilon$  and  $\hat{\Phi}(1)$  is negligible when  $M$  is large enough. Actually, the first integral in (48) can be evaluated as  $\frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}$  because the expression inside the integral is symmetric with respect to  $M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M}$  and  $\Phi(x) + \Phi(-x) = 1$  for any  $x \in \mathbb{R}$ . For large  $M$ , the second integral in (48) can be approximated as:

$$\begin{aligned} &\frac{\sigma}{\mu}\sqrt{M} \int_{Q^{-1}(\epsilon)}^{\frac{\sqrt{M}\mu}{\sigma}} Q(x) dx \stackrel{(a)}{\approx} \frac{\sigma}{\mu}\sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx - \frac{\sigma}{2.5\mu}\sqrt{M} \int_{\frac{\sqrt{M}\mu}{\sigma}}^{\infty} \frac{e^{-\frac{x^2}{2}}}{x} dx \\ &= \frac{\sigma}{\mu}\sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx - \frac{\sigma}{5\mu}\sqrt{M} E_1\left(\frac{M\mu^2}{2\sigma^2}\right) \approx \frac{\sigma}{\mu}\sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x) dx \end{aligned} \quad (49)$$

where (a) follows from [36]: when  $x$  is positive and large enough,  $Q(x) \approx \frac{e^{-\frac{x^2}{2}}}{2.5x}$ . The last line holds because  $\frac{\sigma}{5\mu}\sqrt{M}E_1\left(\frac{M\mu^2}{2\sigma^2}\right) \approx 0$  when  $M$  is sufficiently large [37]. This finally yields:

$$\mathbb{E}[X] \approx \mathbb{E}[\tilde{X}] \approx M - \frac{\sigma}{\mu}Q^{-1}(\epsilon)\sqrt{M} - \frac{\sigma}{\mu}\sqrt{M} \int_{Q^{-1}(\epsilon)}^{\infty} Q(x)dx + 0.5(1 - \epsilon). \quad (50)$$

## APPENDIX C

### PROOF OF THEOREM 3

In order to prove the theorem, we first establish the following lemma:

*Lemma 1:* If the initial rate  $R_{\text{init}}$  has a pre-log of  $r$ , i.e.,  $\lim_{\text{SNR} \rightarrow \infty} \frac{R_{\text{init}}}{\log_2 \text{SNR}} = r$ , then

$$\lim_{\text{SNR} \rightarrow \infty} \mathbb{P} \left[ \sum_{i=1}^k \log_2(1 + \text{SNR}|h_i|^2) \leq R_{\text{init}} \right] = \begin{cases} 1, & \text{for } k < r \text{ and } k \in \mathbb{Z}^+ \\ 0, & \text{for } k > r \text{ and } k \in \mathbb{Z}^+ \end{cases} \quad (51a)$$

*Proof:* For notational convenience we use  $\gamma_i$  to denote the quantity  $\text{SNR}|h_i|^2$ . We prove the first result by using the fact that  $\sum_{i=1}^k \log_2(1 + \gamma_i) \leq k \log_2(1 + \max_{i=1, \dots, k} \gamma_i)$  which yields:

$$\begin{aligned} \mathbb{P} \left[ \sum_{i=1}^k \log_2(1 + \gamma_i) \leq R_{\text{init}} \right] &\geq \mathbb{P} \left[ k \log_2 \left( 1 + \max_{i=1, \dots, k} \gamma_i \right) \leq R_{\text{init}} \right] \\ &= \left( 1 - e^{-\frac{2^{\frac{R_{\text{init}}}{k}} - 1}{\text{SNR}}} \right)^k = \left( 1 - e^{-2^{\frac{R_{\text{init}}}{k} - \log_2(\text{SNR})} + \frac{1}{\text{SNR}}} \right)^k, \end{aligned} \quad (52)$$

where the first equality follows because the  $\gamma_i$ 's are i.i.d. exponential with mean  $\text{SNR}$ . The exponent  $\frac{R_{\text{init}}}{k} - \log_2(\text{SNR})$  behaves as  $\frac{r-k}{k} \log_2(\text{SNR})$ . If  $k < r$  this exponent goes to infinity. Because the  $\frac{1}{\text{SNR}}$  term vanishes,  $e$  is raised to a power converging to  $-\infty$ , and thus (52) converges to 1. This yields the result in (51a). To prove (51b) we combine the property  $\sum_{i=1}^k \log_2(1 + \gamma_i) \geq k \log_2(1 + \min_{i=1, \dots, k} \gamma_i)$  with the same argument as above:

$$\mathbb{P} \left[ \sum_{i=1}^k \log_2(1 + \gamma_i) \leq R_{\text{init}} \right] \leq \mathbb{P} \left[ k \log_2 \left( 1 + \min_{i=1, \dots, k} \gamma_i \right) \leq R_{\text{init}} \right] = 1 - e^{-k \left( 2^{\frac{R_{\text{init}}}{k} - \log_2(\text{SNR})} - \frac{1}{\text{SNR}} \right)}$$

If  $k > r$ ,  $e$  is raised to a power that converges to 0 and thus we get (51b).  $\blacksquare$

We now move on to the proof of the theorem. Using the expression for  $\mathbb{E}[X]$  in (21) we have:

$$\lim_{\text{SNR} \rightarrow \infty} \mathbb{E}[X] = \lim_{\text{SNR} \rightarrow \infty} 1 + \mathbb{P}[\log_2(1 + \gamma_1) \leq R_{\text{init}}] + \dots + \mathbb{P} \left[ \sum_{i=1}^{L-1} \log_2(1 + \gamma_i) \leq R_{\text{init}} \right]. \quad (53)$$

Because  $R_{\text{init}} = MC_\epsilon^M$  has a pre-log of  $M$ , the lemma implies that each of the terms converge to one and thus  $\lim_{\text{SNR} \rightarrow \infty} \mathbb{E}[X] = M$ .

In terms of the high-SNR offset we have:

$$\begin{aligned}
\lim_{\text{SNR} \rightarrow \infty} [C_\epsilon^{\text{IR},M}(\text{SNR}) - C_\epsilon^M(\text{SNR})] &= \lim_{\text{SNR} \rightarrow \infty} \left[ \frac{R_{\text{init}}}{C_\epsilon^M(\text{SNR})} - \mathbb{E}[X] \right] \frac{C_\epsilon^M(\text{SNR})}{\mathbb{E}[X]} \\
&\stackrel{(a)}{\leq} \lim_{\text{SNR} \rightarrow \infty} [M - \mathbb{E}[X]] C_\epsilon^M(\text{SNR}) \\
&= \lim_{\text{SNR} \rightarrow \infty} [M - \mathbb{E}[X]] (\log_2(\text{SNR}) + O(1)) \\
&\stackrel{(b)}{=} \lim_{\text{SNR} \rightarrow \infty} [M - \mathbb{E}[X]] \log_2(\text{SNR}), \tag{54}
\end{aligned}$$

where (a) holds because  $\mathbb{E}[X] \geq 1$  and  $R_{\text{init}} = MC_\epsilon^M(\text{SNR})$ , and (b) holds because  $\mathbb{E}[X] \rightarrow M$  and therefore the  $O(1)$  term does not effect the limit.

Because the additive terms defining  $\mathbb{E}[X]$  in (21) are decreasing, we lower bound  $\mathbb{E}[X]$  as

$$\mathbb{E}[X] \geq M\mathbb{P} \left[ \sum_{i=1}^{M-1} \log_2(1 + \gamma_i) \leq R_{\text{init}} \right] \geq M \left( 1 - e^{-2^{\frac{R_{\text{init}}}{M-1} - \log_2(\text{SNR})} + \frac{1}{\text{SNR}}} \right)^{M-1}, \tag{55}$$

where the last inequality follows from (52). Plugging this bound into (54) yields:

$$\begin{aligned}
\lim_{\text{SNR} \rightarrow \infty} [M - \mathbb{E}[X]] \log_2(\text{SNR}) &\leq \lim_{\text{SNR} \rightarrow \infty} M \left( 1 - \left( 1 - e^{-2^{\frac{R_{\text{init}}}{M-1} - \log_2(\text{SNR})} + \frac{1}{\text{SNR}}} \right)^{M-1} \right) \log_2(\text{SNR}) \\
&= \lim_{\text{SNR} \rightarrow \infty} -M \log_2(\text{SNR}) \sum_{j=1}^{M-1} \binom{M-1}{j} \left( -e^{-2^{\frac{R_{\text{init}}}{M-1} - \log_2(\text{SNR})} + \frac{1}{\text{SNR}}} \right)^j
\end{aligned}$$

where the last line follows from the binomial expansion. Because  $R_{\text{init}}$  has a pre-log of  $M$ , each of the terms is of the form  $\alpha \log_2(\text{SNR}) e^{-\text{SNR}^\beta}$  for some  $\beta > 0$  and some constant  $\alpha$ , and thus the RHS of the last line is zero. Because  $C_\epsilon^{\text{IR},M}(\text{SNR}) \geq C_\epsilon^M(\text{SNR})$ , this shows  $\lim_{\text{SNR} \rightarrow \infty} [C_\epsilon^{\text{IR},M}(\text{SNR}) - C_\epsilon^M(\text{SNR})] = 0$ .

#### APPENDIX D

#### PROOF OF THEOREM 4

In order to prove that rate optimization does not increase the high-SNR offset, we need to consider all possible choices of the initial rate  $R_{\text{init}}$ . We begin by considering all choices of  $R_{\text{init}}$  with a pre-log of  $M$ , i.e., satisfying  $\lim_{\text{SNR} \rightarrow \infty} \frac{R_{\text{init}}}{\log_2 \text{SNR}} = M$ . Because the proof of convergence of  $\mathbb{E}[X]$  in the proof of Theorem 3 only requires  $R_{\text{init}}$  to have a pre-log of  $M$ , we have  $\mathbb{E}[X] \rightarrow M$ . To bound the offset, we write the rate as  $R_{\text{init}} = MC_\epsilon^M(\text{SNR}) - f(\text{SNR})$  where  $f(\text{SNR})$  is strictly positive and sub-logarithmic (because the pre-log is  $M$ ), and thus the rate offset is:

$$\frac{MC_\epsilon^M(\text{SNR}) - f(\text{SNR})}{\mathbb{E}[X]} - C_\epsilon^M(\text{SNR}) = (M - \mathbb{E}[X]) \frac{C_\epsilon^M(\text{SNR})}{\mathbb{E}[X]} - \frac{f(\text{SNR})}{\mathbb{E}[X]}. \tag{56}$$

By the same argument as in Appendix C, the first term is upper bounded by zero in the limit. Therefore:

$$\lim_{\text{SNR} \rightarrow \infty} \frac{R_{\text{init}}}{\mathbb{E}[X]} - C_\epsilon^M(\text{SNR}) \leq \lim_{\text{SNR} \rightarrow \infty} -\frac{f(\text{SNR})}{\mathbb{E}[X]}. \quad (57)$$

Relative to  $C_\epsilon^M(\text{SNR})$ , the offset is either strictly negative (if  $f(\text{SNR})$  is bounded) or goes to negative infinity. In either case a strictly worse offset is achieved.

Let us now consider pre-log factors, denoted by  $r$ , strictly smaller than  $M$  (i.e.,  $r < M$ ). We first consider non-integer values of  $r$ . By Lemma 1 the first  $1 + \lfloor r \rfloor$  terms in the expression for  $\mathbb{E}[X]$  converge to one while the other terms go to zero. Therefore  $\mathbb{E}[X] \rightarrow 1 + \lfloor r \rfloor = \lfloor r \rfloor$ . The long-term transmitted rate, given by  $\frac{R_{\text{init}}}{\mathbb{E}[X]}$ , therefore has pre-log equal to  $\frac{r}{\lfloor r \rfloor}$ . This quantity is strictly smaller than one, and therefore a non-integer  $r$  yields average rate with a strictly suboptimal pre-log factor.

We finally consider integer values of  $r$  satisfying  $r < M$ . In this case we must separately consider rates of the form  $R_{\text{init}} = r \log \text{SNR} \pm O(1)$  versus those of the form  $R_{\text{init}} = r \log \text{SNR} \pm o(\log \text{SNR})$ . Here we use  $o(\log \text{SNR})$  to denote terms that are sub-logarithmic and that go to positive infinity; note that we also explicitly denote the sign of the  $O(1)$  or  $o(\log \text{SNR})$  terms. We first consider  $R_{\text{init}} = r \log \text{SNR} \pm O(1)$ . By Lemma 1 the terms corresponding to  $k = 0, \dots, r - 1$  in the expression for  $\mathbb{E}[X]$  converge to one, while the terms corresponding to  $k = r + 1, \dots, M - 1$  go to one. Furthermore, the term corresponding to  $k = r$  converges to a strictly positive constant denoted  $\delta$ :

$$\begin{aligned} \delta &= \lim_{\text{SNR} \rightarrow \infty} \mathbb{P} \left[ \sum_{i=1}^r \log_2(1 + \text{SNR}|h_i|^2) \leq r \log_2(\text{SNR}) \pm O(1) \right] \\ &= \lim_{\text{SNR} \rightarrow \infty} \mathbb{P} \left[ \sum_{i=1}^r \log_2(\text{SNR}|h_i|^2) \leq r \log_2(\text{SNR}) \pm O(1) \right] = \mathbb{P} \left[ \sum_{i=1}^r \log_2(|h_i|^2) \leq \pm O(1) \right] \end{aligned} \quad (58)$$

where the second equality follows from [26].  $\delta$  is strictly positive because the support of  $\log_2(|h_i|^2)$ , and thus of the sum, is the entire real line. As a result,  $\mathbb{E}[X] \rightarrow r + \delta$ , which is strictly larger than  $r$ . The pre-log of the average rate is then  $\frac{r}{r+\delta} < 1$ , and so this choice of initial rate is also sub-optimal.

If  $R_{\text{init}} = r \log \text{SNR} + o(\log \text{SNR})$  the terms in the  $\mathbb{E}[X]$  expression behave largely the same as above except that the  $k = r$  term converges to one because the  $O(1)$  term in (58) is replaced with a quantity tending to positive infinity. Therefore  $\mathbb{E}[X] \rightarrow r + 1$ , which also yields a sub-optimal pre-log of  $\frac{r}{r+1} < 1$ .

We are thus finally left with the choice  $R_{\text{init}} = r \log \text{SNR} - o(\log \text{SNR})$ . This is the same as the above case except that the  $k = r$  term converges to zero. Therefore  $\mathbb{E}[X] \rightarrow r$ , and thus the achieved pre-log is one. In this case we must explicitly consider the rate offset, which is written as:

$$\frac{R_{\text{init}}}{\mathbb{E}[X]} - C_\epsilon^M = \frac{r \log \text{SNR} - o(\log \text{SNR})}{\mathbb{E}[X]} - C_\epsilon^M = \frac{r \log \text{SNR}}{\mathbb{E}[X]} - C_\epsilon^M - \frac{o(\log \text{SNR})}{\mathbb{E}[X]}. \quad (59)$$

Using essentially the same proof as for Theorem 3, the difference between the first two terms is upper bounded by zero in the limit of  $\text{SNR} \rightarrow \infty$ . Thus, the rate offset goes to negative infinity.

Because we have shown each choice of  $R_{\text{init}}$  (except  $R_{\text{init}} = MC_{\epsilon}^M$ ) achieves either a strictly sub-optimal pre-log or the correct pre-log but a strictly negative offset, this proves both parts of the theorem.

## REFERENCES

- [1] D. Chase, "Class of algorithms for decoding block codes with channel measurement information," *IEEE Trans. Inf. Theory*, vol. 18, no. 1, pp. 170–182, 1972.
- [2] M. Sauter, *Communications Systems for the Mobile Information Society*. John Wiley, 2006.
- [3] J. G. Andrews, A. Ghosh, and R. Muhamed, *Fundamentals of Wimax: Understanding Broadband Wireless Networking*, 1st ed. Prentice Hall PTR, 2007.
- [4] L. Favalli, M. Lanati, and P. Savazzi, *Wireless Communications 2007 CNIT Thyrranian Symposium*, 1st ed. Springer, 2007.
- [5] S. Sesia, G. Caire, and G. Vivier, "Incremental redundancy hybrid ARQ schemes based on low-density parity-check codes," *IEEE Trans. Commun.*, vol. 52, no. 8, pp. 1311–1321, Aug. 2004.
- [6] D. Costello, J. Hagenauer, H. Imai, and S. Wicker, "Applications of error-control coding," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2531–2560, Oct. 1998.
- [7] E. Soljanin, N. Varnica, and P. Whiting, "Incremental redundancy hybrid ARQ with LDPC and raptor codes," *submitted to IEEE Trans. Inf. Theory*, 2005.
- [8] C. Lott, O. Milenkovic, and E. Soljanin, "Hybrid ARQ: theory, state of the art and future directions," *IEEE Inf. Theory Workshop*, pp. 1–5, Jul. 2007.
- [9] G. Caire and D. Tuninetti, "The throughput of hybrid-ARQ protocols for the Gaussian collision channel," *IEEE Trans. Inf. Theory*, vol. 47, no. 4, pp. 1971–1988, Jul. 2001.
- [10] H. E. Gamal, G. Caire, and M. E. Damen, "The MIMO ARQ channel: diversity-multiplexing-delay tradeoff," *IEEE Trans. Inf. Theory*, pp. 3601–3621, 2006.
- [11] A. Chuang, A. Guillén i Fàbregas, L. K. Rasmussen, and I. B. Collings, "Optimal throughput-diversity-delay tradeoff in MIMO ARQ block-fading channels," *IEEE Trans. Inf. Theory*, vol. 54, no. 9, pp. 3968–3986, Sep. 2008.
- [12] C. Shen, T. Liu, and M. P. Fitz, "Aggressive transmission with ARQ in quasi-static fading channels," *Proc. of IEEE Int'l Conf. in Commun. (ICC'08)*, pp. 1092–1097, May 2008.
- [13] R. Narasimhan, "Throughput-delay performance of half-duplex hybrid-ARQ relay channels," *Proc. of IEEE Int'l Conf. in Commun. (ICC'08)*, pp. 986–990, May 2008.
- [14] X. Tang, R. Liu, P. Spasojevic, and H. V. Poor, "On the Throughput of Secure Hybrid-ARQ Protocols for Gaussian Block-Fading Channels," *IEEE Trans. Inf. Theory*, vol. 55, no. 4, pp. 1575–1591, Apr. 2009.
- [15] L. Zheng and D. Tse, "Diversity and multiplexing: a fundamental tradeoff in multiple-antenna channels," *IEEE Trans. Inf. Theory*, vol. 49, no. 5, pp. 1073–1096, May 2003.
- [16] A. Lozano and N. Jindal, "Transmit diversity v. spatial multiplexing in modern MIMO systems," *submitted to IEEE Trans. Wireless Commun.*, 2008.
- [17] P. Wu and N. Jindal, "Coding Versus ARQ in Fading Channels: How reliable should the PHY be?" *submitted to Proc. of IEEE Globe Commun. Conf. (Globecom'09)*, 2009.

- [18] T. A. Courtade and R. D. Wesel, "A cross-layer perspective on rateless coding for wireless channels," *to appear at Proc. of IEEE Int'l Conf. in Commun. (ICC'09)*, Jun. 2009.
- [19] D. Tse and P. Viswanath, *Fundamentals of Wireless Communications*. Cambridge University, 2005.
- [20] G. Carie, G. Taricco, and E. Biglieri, "Optimum power control over fading channels," *IEEE Trans. Inf. Theory*, vol. 45, no. 5, pp. 1468–1489, Jul. 1999.
- [21] A. Guillén i Fàbregas and G. Caire, "Coded modulation in the block-fading channel: coding theorems and code construction," *IEEE Trans. Inf. Theory*, vol. 52, no. 1, pp. 91–114, Jan. 2006.
- [22] N. Prasad and M. K. Varanasi, "Outage theorems for MIMO block-fading channels," *IEEE Trans. Inf. Theory*, vol. 52, no. 12, pp. 5284–5296, Dec. 2006.
- [23] E. Malkamäki and H. Leib, "Coded diversity on block-fading channels," *IEEE Trans. Inf. Theory*, vol. 45, no. 2, pp. 771–781, Mar. 1999.
- [24] S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Trans. Inf. Theory*, vol. 40, no. 4, pp. 1147–1157, Jul. 1994.
- [25] S. Shamai and S. Verdú, "The impact of frequency-flat fading on the spectral efficiency of CDMA," *IEEE Trans. Inf. Theory*, vol. 47, no. 5, pp. 1302–1327, May 2001.
- [26] N. Prasad and M. Varanasi, "MIMO outage capacity in the high SNR regime," *Proc. of IEEE Int'l Symp. on Inform. Theory (ISIT'05)*, pp. 656–660, Sep. 2005.
- [27] A. Lozano, A. M. Tulino, and S. Verdú, "High-SNR power offset in multiantenna communication," *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4134–4151, Dec. 2005.
- [28] J. Salo, H. E. Sallabi, and P. Vainikainen, "The distribution of the product of independent Rayleigh random variables," *IEEE Trans. Antennas Propag.*, vol. 54, no. 2, pp. 639–643, Feb. 2006.
- [29] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, 5th ed. Academic, 1994.
- [30] G. Barriac and U. Madhow, "Characterizing outage rates for space-time communication over wideband channels," *IEEE Trans. Commun.*, vol. 52, no. 4, pp. 2198–2207, Dec. 2004.
- [31] P. J. Smith and M. Shañi, "On a Gaussian approximation to the capacity of wireless MIMO systems," *Proc. of IEEE Int'l Conf. in Commun. (ICC'02)*, pp. 406–410, Apr. 2002.
- [32] M. S. Alouini and A. J. Goldsmith, "Capacity of Rayleigh fading channels under different adaptive transmission and diversity-combining techniques," *IEEE Trans. Veh. Technol.*, vol. 48, no. 4, pp. 1165–1181, Jul. 1999.
- [33] M. R. McKay, P. J. Smith, H. A. Suraweera, and I. B. Collings, "On the mutual information distribution of OFDM-based spatial multiplexing: exact variance and outage approximation," *IEEE Trans. Inf. Theory*, vol. 54, no. 7, pp. 3260–3278, 2008.
- [34] B. Fristedt and L. Gray, *A Modern Approach to Probability Theory*. Birkhäuser Boston, 1997.
- [35] M. Meyer, H. Wiemann, M. Sagfors, J. Torsner, and J.-F. Cheng, "ARQ concept for the UMTS long-term evolution," *Proc. of IEEE Vehic. Technol. Conf. (VTC'2006 Fall)*, Sep. 2006.
- [36] N. Kingsbury, "Approximation formula for the Gaussian error integral,  $Q(x)$ ," <http://cnx.org/content/m11067/latest/>.
- [37] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover Publications, 1964.

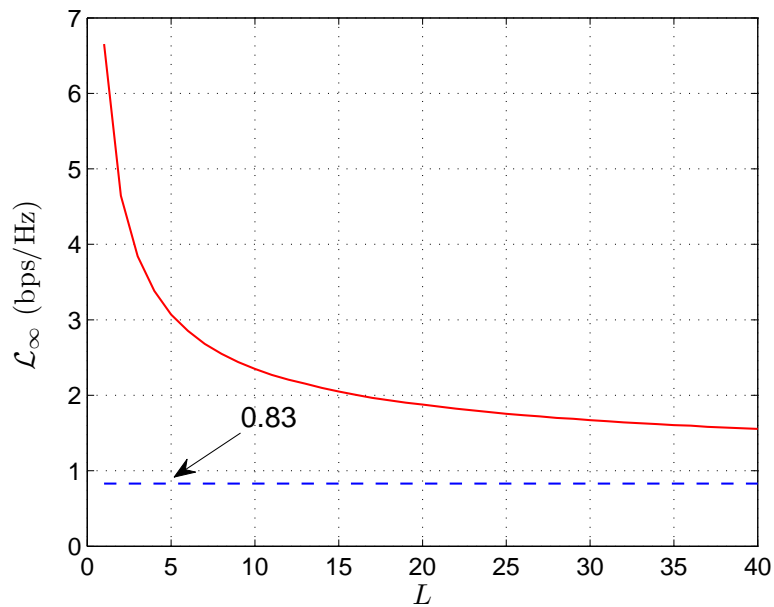


Fig. 1. High SNR rate offset  $\mathcal{L}_\infty$  (bps/Hz) versus diversity order  $L$  for  $\epsilon = 0.01$

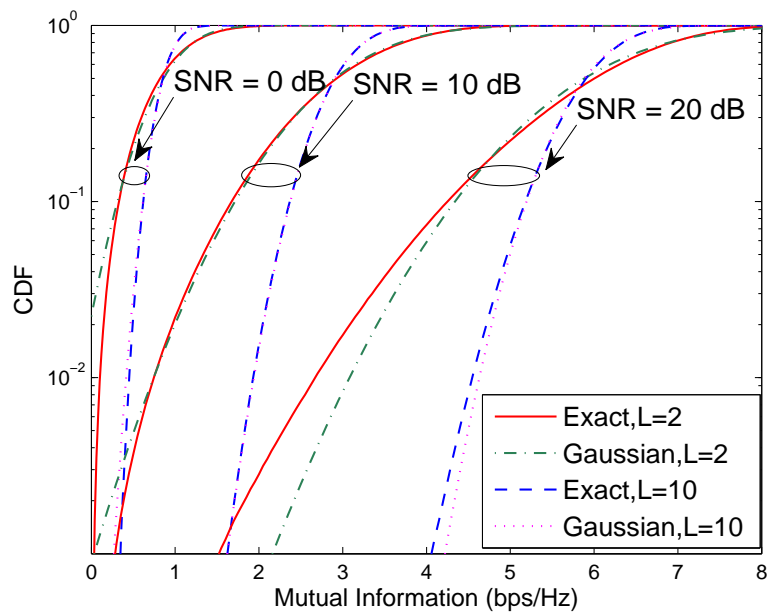


Fig. 2. CDF's of mutual information ( $\frac{1}{2} \sum_{i=1}^2 \log_2(1 + \text{SNR}|h_i|^2)$  and  $\frac{1}{10} \sum_{i=1}^{10} \log_2(1 + \text{SNR}|h_i|^2)$  respectively) for  $L = 2$  and  $L = 10$  at SNR = 0, 10, and 20 dB

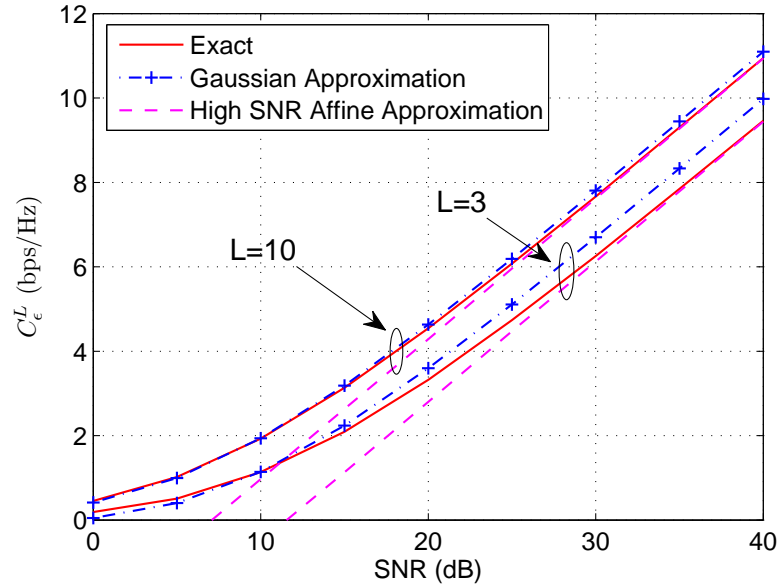


Fig. 3.  $\epsilon$ -outage capacity  $C_\epsilon^L$  (bps/Hz) versus SNR (dB) for  $\epsilon = 0.01$

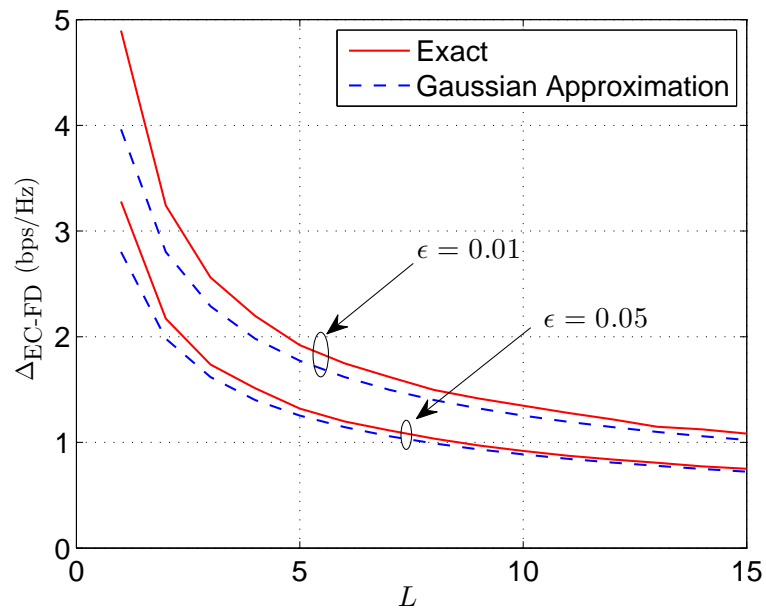


Fig. 4. Ergodic capacity -  $\epsilon$ -outage capacity difference  $\Delta_{\text{EC-FD}}$  (bps/Hz) versus diversity order  $L$ , at SNR = 20 dB

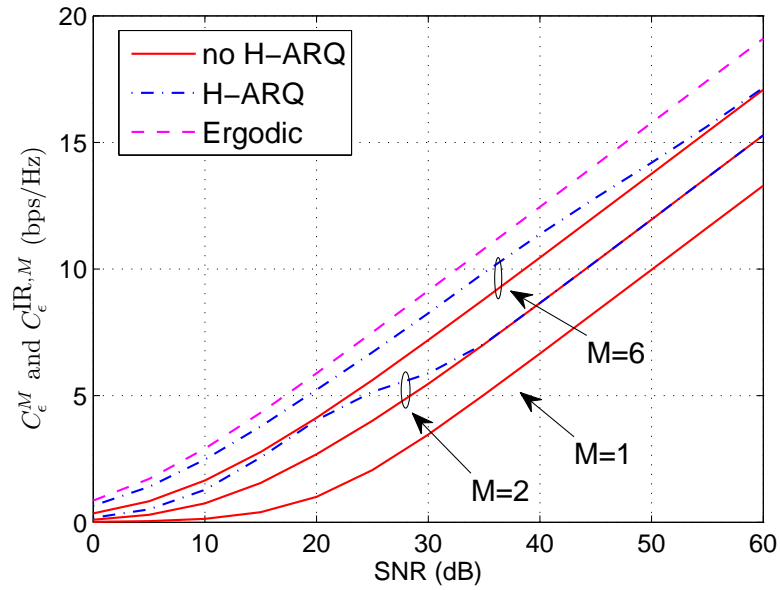


Fig. 5. Ergodic capacity, H-ARQ rate, and non-H-ARQ rate (bps/Hz) versus SNR (dB) for  $\epsilon = 0.01$

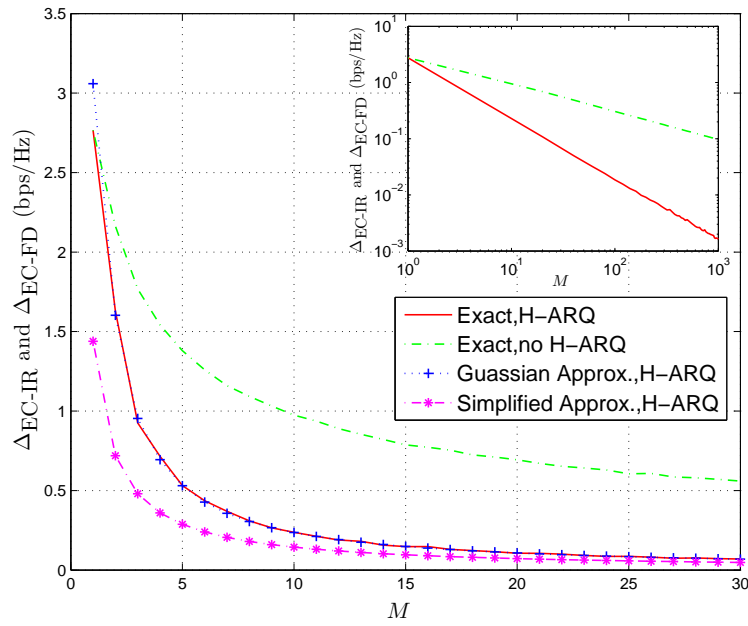


Fig. 6. Ergodic capacity - H-ARQ rate difference  $\Delta_{EC-IR}$  (bps/Hz) and ergodic capacity - non-H-ARQ rate difference  $\Delta_{EC-FD}$  (bps/Hz) versus  $M$  for  $\epsilon = 0.01$  at SNR = 10 dB

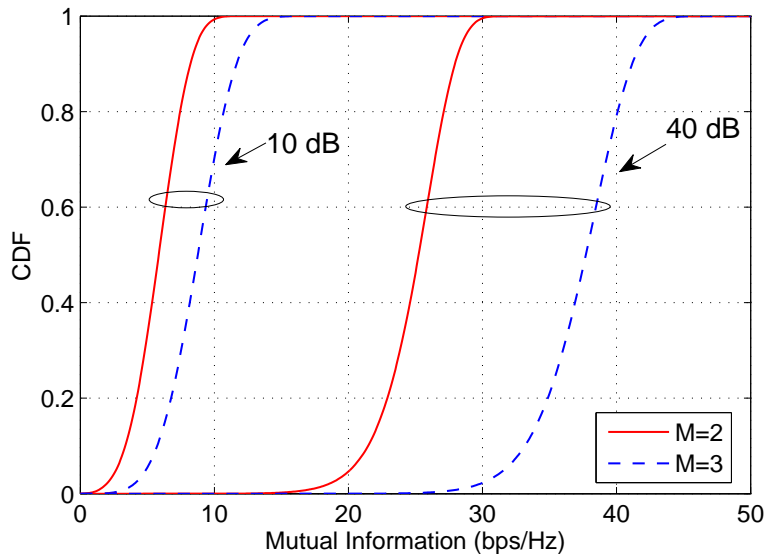


Fig. 7. CDF's of accumulated mutual information over 2 and 3 rounds (random variables  $\sum_{i=1}^2 \log_2(1 + \text{SNR}|h_i|^2)$  and  $\sum_{i=1}^3 \log_2(1 + \text{SNR}|h_i|^2)$ , respectively) at SNR = 10 and 40 dB

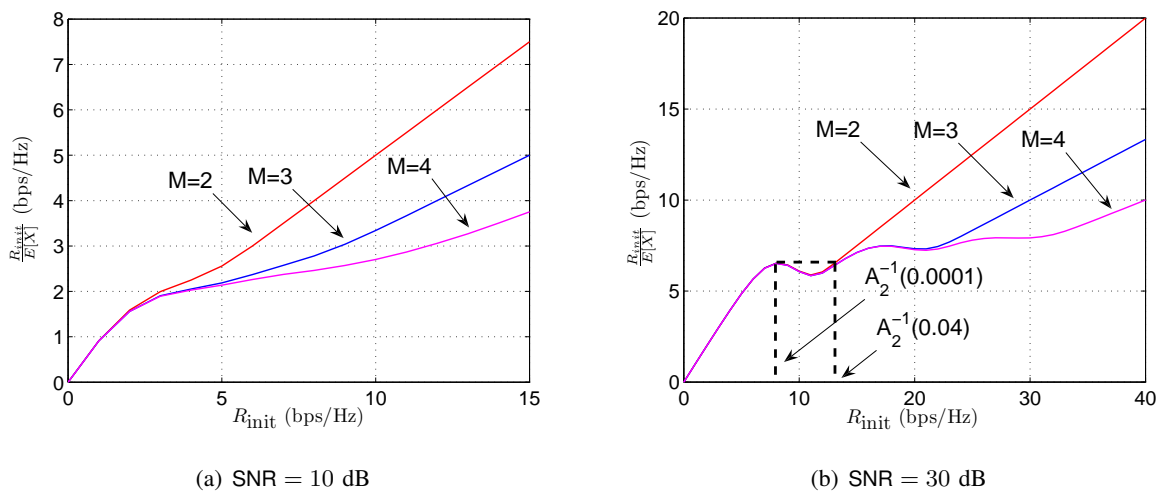


Fig. 8. H-ARQ rate  $\frac{R_{\text{init}}}{\mathbb{E}[X]}$  (bps/Hz) versus initial rate  $R_{\text{init}}$  (bps/Hz)

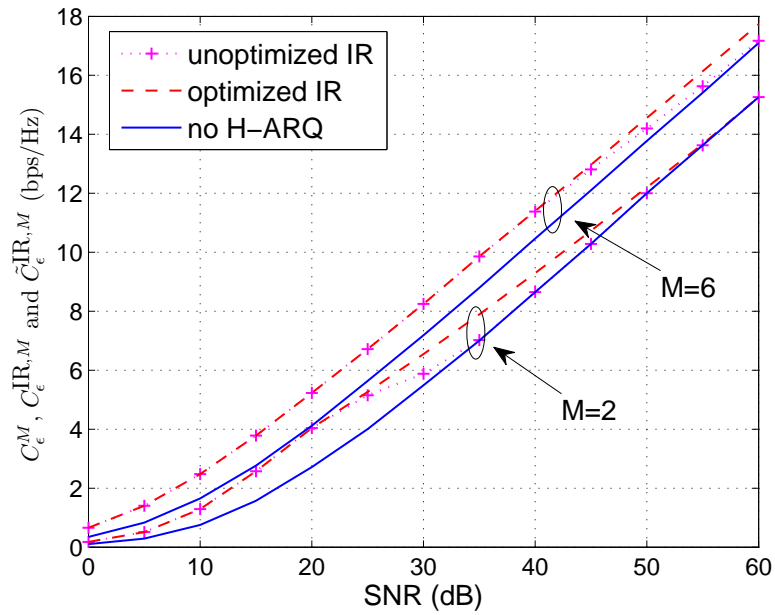


Fig. 9. H-ARQ rate with/without an optimized initial rate (bps/Hz) and non-H-ARQ rate (bps/Hz) versus SNR (dB) for  $\epsilon = 0.01$

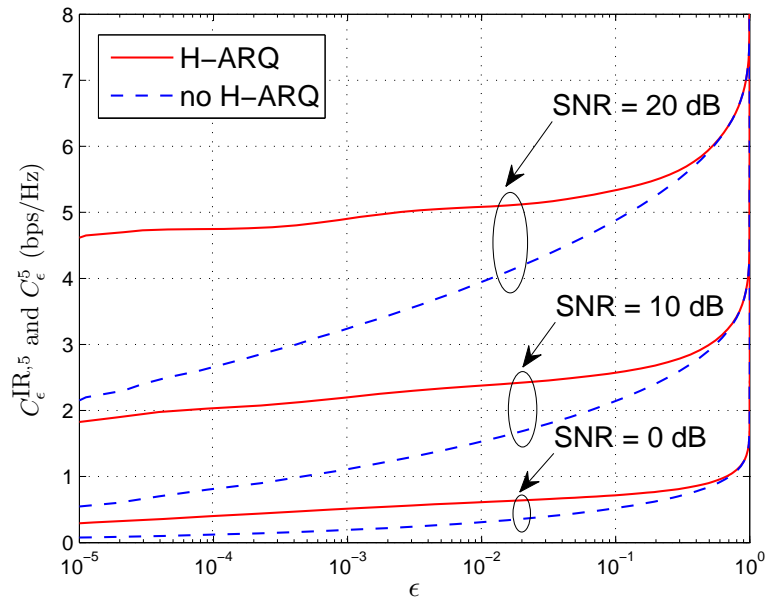


Fig. 10. H-ARQ rate (bps/Hz) and non-H-ARQ rate (bps/Hz) versus outage probability  $\epsilon$  for  $M = 5$  at SNR = 0 and 10 and 20 dB

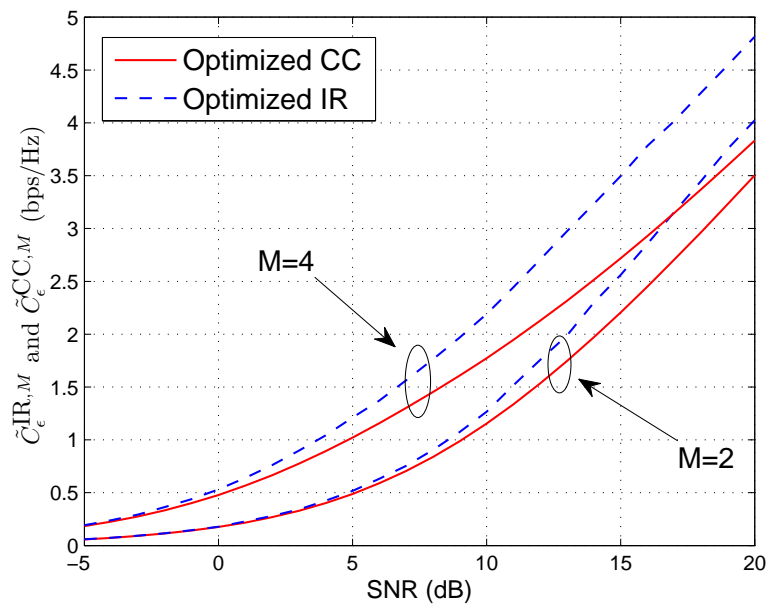


Fig. 11. Optimized H-ARQ rate  $\tilde{C}_\epsilon^{\text{IR},M}$  (bps/Hz) and  $\tilde{C}_\epsilon^{\text{CC},M}$  (bps/Hz) versus SNR (dB) for  $\epsilon = 0.01$