

# Channel Coding Rate in the Finite Blocklength Regime

Yury Polyanskiy, *Student Member, IEEE*, H. Vincent Poor, *Fellow, IEEE*, and Sergio Verdú, *Fellow, IEEE*

**Abstract**—This paper investigates the maximal channel coding rate achievable at a given blocklength and error probability. For general classes of channels new achievability and converse bounds are given, which are tighter than existing bounds for wide ranges of parameters of interest, and lead to tight approximations of the maximal achievable rate for blocklengths  $n$  as short as 100. It is also shown analytically that the maximal rate achievable with error probability  $\epsilon$  is closely approximated by  $C - \sqrt{\frac{V}{n}} Q^{-1}(\epsilon)$  where  $C$  is the capacity,  $V$  is a characteristic of the channel referred to as *channel dispersion*, and  $Q$  is the complementary Gaussian cumulative distribution function.

**Index Terms**—Achievability, channel capacity, coding for noisy channels, converse, finite blocklength regime, Shannon theory.

## I. INTRODUCTION

THE proof of the channel coding theorem involves three stages:

- *Converse*: an upper bound on the size of any code with given arbitrary blocklength and error probability.
- *Achievability*: a lower bound on the size of a code that can be guaranteed to exist with given arbitrary blocklength and error probability.
- *Asymptotics*: the bounds on the log size of the code normalized by blocklength asymptotically coincide as a result of the law of large numbers (memoryless channels) or another ergodic theorem (for channels with memory).

As propounded in [1], it is pedagogically sound to separate clearly the third stage from the derivation of the upper and lower bounds:

- The bounds need not impose assumptions on the channel such as memorylessness, stationarity, and ergodicity.
- The key information theoretic arguments are used mainly in the converse and achievability bounds.
- The bounds can be extremely useful in assessing the highest rate that can be achieved when operating with a given blocklength and error probability.

The strong form of the coding theorem establishes that for a general class of channels that behave ergodically [2], the

channel capacity is the largest rate at which information can be transmitted regardless of the desired error probability, provided that the blocklength is allowed to grow without bound. In practice, it is of vital interest to assess the backoff from capacity required to sustain the desired error probability at a given fixed finite blocklength. Unfortunately, no guidance to answer that question is offered either by the strong version of the coding theorem, or by the reliability function, which gives the asymptotic exponential decay of error probability when transmitting at any given fraction of capacity.

In the nonasymptotic regime, there are no exact formulas for the maximal rate sustainable as a function of blocklength and error probability. In this paper, we show several new achievability and converse bounds which bound the fundamental limits tightly for blocklengths as short as 100. Together with normal approximations, the bounds also show that in the finite blocklength regime, the backoff from channel capacity  $C$  is accurately and succinctly characterized by a parameter that we refer to as the *channel dispersion*  $V$ , which measures the stochastic variability of the channel relative to a deterministic channel with the same capacity. Specifically, the finite blocklength coding rate is approximated by<sup>1</sup>

$$\frac{1}{n} \log M^*(n, \epsilon) \approx C - \sqrt{\frac{V}{n}} Q^{-1}(\epsilon) \quad (1)$$

where  $n$  is the blocklength and  $\epsilon$  is the error probability.

Since Shannon established the convergence of optimal coding rate to capacity, there has been some work devoted to the assessment of the penalty incurred by finite blocklength. Foremost, Shannon [3] provided tight bounds for the additive white Gaussian noise (AWGN) channel that were studied numerically by Slepian [4] (cf. also [5] and [6]). Recently, with the advent of sparse-graph codes, a number of works [8]–[11] have studied the SNR penalty as a function of blocklength in order to improve the assessment of the suboptimality of a given code with respect to the fundamental limit at that particular blocklength rather than the asymptotic limit embodied in the channel capacity. Approximations of the type in (1) have been studied in [7], [23], [29]–[31].

The major existing achievability and converse bounds are reviewed in Section II along with refined asymptotic expansions of achievable rate. Section III gives our new lower and upper bounds on the maximal rate achievable for a given blocklength and error probability. The lower bounds are based on three different constructive approaches that lead, respectively, to the RCU (random-coding union) bound, the DT (dependency testing) bound, and the  $\kappa/\beta$  bound based on the Neyman–Pearson

Manuscript received November 14, 2008; revised October 22, 2009. Current version published April 21, 2010. This work was supported in part by the National Science Foundation by Grants CCF-06-35154, CCF-07-28445, and CNS-09-05398.

The authors are with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: ypolyans@princeton.edu; poor@princeton.edu; verdu@princeton.edu).

Communicated by G. Kramer, Associate Editor for Shannon Theory.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIT.2010.2043769

<sup>1</sup>As usual,  $Q(x) = \int_x^\infty \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$ .

lemma that uses an auxiliary output distribution. Unlike existing achievability bounds, the RCU and DT bounds contain no parameters (other than the input distribution) to be optimized. A general converse upper bound is given as a result of the solution of a minimax problem on the set of input/output distributions. Section IV studies the normal approximation to the maximal achievable rate for discrete memoryless channels and for the additive white Gaussian noise channel, and shows that (1) holds up to a term of  $O((\log n)/n)$  except in rare cases. Throughout Sections III and IV, particular attention is given to the binary erasure channel (BEC), the binary symmetric channel (BSC), and the AWGN channel. Several coding schemes used in practice are compared against the nonasymptotic fundamental limits. The use of the normal approximation as a design tool is illustrated in the context of the optimization of the maximal throughput of a simple automatic repeat request (ARQ) retransmission strategy. Section V summarizes our main findings.

## II. PREVIOUS WORK

Let us consider input and output sets  $A$  and  $B$  and a conditional probability measure  $P_{Y|X} : A \mapsto B$ . We denote a codebook with  $M$  codewords by  $(c_1, \dots, c_M) \in A^M$ . A (possibly randomized) decoder is a random transformation  $P_{Z|Y} : B \mapsto \{0, 1, \dots, M\}$  (where '0' indicates that the decoder chooses "error"). A codebook with  $M$  codewords and a decoder that satisfies  $P_{Z|X}(m|c_m) \geq 1 - \epsilon$  for  $m \in \{1, \dots, M\}$  are called an  $(M, \epsilon)$ -code (maximal probability of error). If the messages are equiprobable, the average error probability is

$$1 - \frac{1}{M} \sum_{m=1}^M P_{Z|X}(m|c_m).$$

A codebook and a decoder whose average probability of error is smaller than  $\epsilon$  are called an  $(M, \epsilon)$ -code (average probability of error). In the application of our results, we will take  $A$  and  $B$  to be  $n$ -fold Cartesian products of alphabets  $\mathcal{A}$  and  $\mathcal{B}$ , and a channel to be a sequence of conditional probabilities  $\{P_{Y^n|X^n} : \mathcal{A}^n \rightarrow \mathcal{B}^n\}$  [2]. An  $(M, \epsilon)$  code for  $\{\mathcal{A}^n, \mathcal{B}^n, P_{Y^n|X^n}\}$  is called an  $(n, M, \epsilon)$  code. The maximal code size achievable with a given error probability and blocklength is denoted by

$$M^*(n, \epsilon) = \max\{M : \exists \text{ an } (n, M, \epsilon)\text{-code}\}. \quad (2)$$

For the statement and proof of the achievability and converse bounds, it is preferable not to assume that  $A$  and  $B$  have any structure such as a Cartesian product. This has the advantage of avoiding the notational clutter that results from explicitly showing the dimension  $(n)$  of the random variables taking values on  $A$  and  $B$ .

### A. Achievability Bounds Without Codeword Constraints

For a joint distribution  $P_{XY}$  on  $A \times B$  we denote the information density by

$$i(x; y) = \log \frac{dP_{XY}}{d(P_X \times P_Y)}(x, y) \quad (3)$$

$$= \log \frac{dP_{Y|X=x}}{dP_Y}(y) \quad (4)$$

with the understanding that if  $P_{Y|X=x}$  is not absolutely continuous with respect to  $P_Y$  we define  $i(x; y) = +\infty$  for all  $y$  in the singular set, and we define  $i(x; y) = -\infty$  for any  $y$  such that  $\frac{dP_{Y|X=x}}{dP_Y} = 0$ .

Feinstein's [13] achievability bound for maximal probability of error is given as follows.

*Theorem 1 (Feinstein):* For any distribution  $P_X$ , and any  $\beta > 0$ , there exists an  $(M, \epsilon)$  code (maximal probability of error) such that<sup>2</sup>

$$M \geq \beta(\epsilon - \mathbb{P}[i(X; Y) \leq \log \beta]). \quad (5)$$

Alternatively, Shannon's achievability bound [14] is given as follows.

*Theorem 2 (Shannon):* For any distribution  $P_X$ , and any  $\beta > 0$ , there exists an  $(M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq \frac{M-1}{\beta} + \mathbb{P}[i(X; Y) \leq \log \beta]. \quad (6)$$

It is easy to verify that Theorem 1 implies a slightly weakened version of Theorem 2 where  $M-1$  is replaced by  $M$ ; conversely, Theorem 2 implies the weakened version of Theorem 1 where maximal is replaced by average error probability.

The following upper bound is a reformulation of Gallager's random coding bound [15], in terms of information density.

*Theorem 3 (Gallager):* For any  $P_X$  and  $\lambda \in [0, 1]$ , there exists an  $(M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq M^\lambda \mathbb{E} \left[ \left( \mathbb{E} \left[ \exp \frac{i(\bar{X}; Y)}{1+\lambda} \mid Y \right] \right)^{1+\lambda} \right] \quad (7)$$

where  $P_{XY\bar{X}}(a, b, c) = P_X(a)P_{Y|X=a}(b)P_X(c)$ .

For a memoryless channel (7) turns, after optimization over  $\lambda$ , into

$$\epsilon \leq \exp\{-nE_r(R)\} \quad (8)$$

where  $R = \frac{\log M}{n}$  and  $E_r(R)$  is Gallager's random coding exponent [16].

### B. Achievability Bounds With Linear Codes

For a linear code over the BSC, Poltyrev [17] proved the following upper bound on the probability of error.

*Theorem 4 (Poltyrev):* The maximal probability of error  $P_e$  under maximum likelihood decoding of a linear code<sup>3</sup> with weight distribution<sup>4</sup>  $\{A_w, w = 0, \dots, n\}$  over the BSC with crossover probability  $\delta$  satisfies

$$P_e \leq \sum_{\ell=0}^n \delta^\ell (1-\delta)^{n-\ell} \min \left\{ \binom{n}{\ell}, \sum_{w=0}^n A_w B(\ell, w, n) \right\} \quad (9)$$

<sup>2</sup> $P_X, Q$  and  $P_{Y|X=x}$  denote distributions, whereas  $\mathbb{P}$  is reserved for the probability of an event on the underlying probability space.

<sup>3</sup>At the expense of replacing maximal probability of error by average, the same bound can be shown for a nonlinear code by generalizing the notion of weight distribution.

<sup>4</sup>We define  $A_0$  to be the number of 0-weight codewords in the codebook minus 1. In particular, for a linear codebook with no repeated codewords  $A_0 = 0$ .

where

$$B(\ell, w, n) = \sum_{w/2 \leq t \leq \min\{\ell, w\}} \binom{w}{t} \binom{n-w}{\ell-t}. \quad (10)$$

A  $[k, n]$  linear code is generated by a  $k \times n$  binary matrix. We can average (9) over an equiprobable ensemble of such matrices. Applying Jensen's inequality to pass expectation inside the minimum and noticing that  $\mathbb{E}[A_w] = 2^{k-n} \binom{n}{w}$  we obtain the following achievability bound.

*Theorem 5:* For a BSC with crossover probability  $\delta$  there exists a  $[k, n]$  linear code such that a maximum likelihood decoder has a maximal probability of error  $P_e$  satisfying

$$P_e \leq \sum_{\ell=0}^n \delta^\ell (1-\delta)^{n-\ell} \times \min \left\{ \binom{n}{\ell}, \sum_{w=0}^n 2^{k-n} \binom{n}{w} B(\ell, w, n) \right\} \quad (11)$$

where  $B(\ell, w, n)$  is given by (9).

A negligible improvement to (11) is possible if we average (9) over an ensemble of all full-rank binary matrices instead. Another modification by expurgating low-weight codewords [18] leads to a tightening of (11) when the rate is much lower than capacity.

For the BEC the results of [19, Th. 9] can be used to compute the exact value of the probability of error over an ensemble of all linear codes generated by full-rank  $k \times n$  binary matrices [20].

*Theorem 6 (Ashikhmin):* Given a BEC with erasure probability  $\delta$ , the average probability of error over all binary  $k \times n$  linear codes with full-rank generating matrices chosen equiprobably is equal to (12), shown at the bottom of the page, where

$$\begin{bmatrix} a \\ r \end{bmatrix} \triangleq \prod_{j=0}^{r-1} \frac{2^a - 2^j}{2^r - 2^j}$$

is the number of  $r$ -dimensional subspaces of  $\mathbb{F}_2^a$ .

### C. Achievability Bounds With Codeword Constraints

Suppose that all codewords are required to belong to some set  $F \subset A$ . For example, there might be a cost  $c(x)$  associated with using a particular input vector  $x$ , in which case the set  $F$  might be chosen as

$$F = \{x : c(x) \leq P\}. \quad (13)$$

A cost-constrained generalization of (5) due to Thomasian [21] (see also [22]) in which all the codewords are constrained to belong to  $F$  is

$$M \geq \beta (\epsilon - \mathbb{P}[i(X; Y) \leq \log \beta] - P_X[F^c]). \quad (14)$$

A cost-constrained version of (6) is

$$\epsilon \leq \frac{M-1}{\beta} + \mathbb{P}[i(X; Y) \leq \log \beta] + P_X[F^c]. \quad (15)$$

It should be noted that in both (14), and (15), the auxiliary distribution  $P_X$  is not constrained to take values on  $F$ . Theorem 3 admits the following generalization to the setting with cost constraints.

*Theorem 7 (Gallager, With Cost):* Suppose  $P_X$  is such that

$$\mathbb{E}[c(X)] \leq P \quad (16)$$

and denote

$$\mu(\delta) \triangleq \mathbb{P}[P - \delta \leq c(X) \leq P]. \quad (17)$$

Then, for any  $\delta \in [0, P]$  such that  $\mu(\delta) > 0$ ,  $\lambda \in [0, 1]$  and  $r \geq 0$  there exists an  $(M, \epsilon)$ -code (average probability of error) with codewords in  $F$  given by (13) and such that we have (18), shown at the bottom of the page, where  $P_{XY\bar{X}}(a, b, c) = P_X(a)P_{Y|X=a}(b)P_X(c)$ .

### D. Converse Results

The simplest upper bound on the size of a code as a function of the average error probability follows from Fano's inequality:

*Theorem 8:* Every  $(M, \epsilon)$ -code (average probability of error) for a random transformation  $P_{Y|X}$  satisfies

$$\log M \leq \frac{1}{1-\epsilon} \sup_X I(X; Y) + \frac{1}{1-\epsilon} h(\epsilon) \quad (19)$$

$$P_e = \sum_{i=0}^n \binom{n}{i} \delta^{n-i} (1-\delta)^i \sum_{r=\max\{0, k-n+i\}}^{\min\{k, i\}} \begin{bmatrix} i \\ r \end{bmatrix} \begin{bmatrix} n-i \\ k-r \end{bmatrix} \begin{bmatrix} n \\ k \end{bmatrix}^{-1} 2^{r(n-i-k+r)} (1-2^{r-k}) \quad (12)$$

$$\epsilon \leq M^\lambda \left( \frac{\exp(r\delta)}{\mu(\delta)} \right)^{1+\lambda} \mathbb{E} \left[ \left( \mathbb{E} \left[ \exp \left\{ \frac{i(\bar{X}; Y)}{1+\lambda} + (c(\bar{X}) - P)r \right\} \middle| Y \right] \right)^{1+\lambda} \right] \quad (18)$$

where  $h(x) = -x \log x - (1-x) \log(1-x)$  is the binary entropy function.

A significant improvement under the maximal error probability formalism is supplied by the bound due to Wolfowitz [23].

*Theorem 9 (Wolfowitz):* Every  $(M, \epsilon)$ -code (maximal probability of error) must satisfy

$$M \leq \inf_{\beta > 0} \beta \left( \inf_{x \in \mathcal{A}} P_{Y|X=x} [i(x; Y) < \log \beta] - \epsilon \right)^{-1} \quad (20)$$

provided that the right-hand side (RHS) is not less than 1.

As shown in [24, Th. 7.8.1], this bound leads to the strong converse theorem for the discrete memoryless channel (DMC), even assuming noiseless feedback, namely

$$\log M^*(n, \epsilon) \leq nC + o(n) \quad \forall \epsilon \in (0, 1). \quad (21)$$

Theorem 9 can be further tightened by maximizing the probability therein with respect to the choice of the unconditional output distribution in the definition of information density [25].

The following corollary to Theorem 9 gives another converse bound which also leads to (21), but is too coarse for the purposes of analyzing the fundamental limits in the finite block-length regime.

*Theorem 10 ([16, Th. 5.8.5]):* For an arbitrary discrete memoryless channel of capacity  $C$  and any  $(n, \exp\{nR\}, \epsilon)$  code with rate  $R > C$ , we have

$$\epsilon \geq 1 - \frac{4A}{n(R-C)^2} - \exp\left\{-\frac{n(R-C)}{2}\right\} \quad (22)$$

where  $A > 0$  is constant independent of  $n$  or  $R$ .

The dual of the Shannon–Feinstein bounds in Theorems 1 and 2 (in the unconstrained setting) is given in [2].

*Theorem 11 (Verdú–Han):* Every  $(M, \epsilon)$ -code (average error probability) satisfies

$$\epsilon \geq \sup_{\beta > 0} \left\{ \inf_{P_X} \mathbb{P} [i(X; Y) \leq \log \beta] - \frac{\beta}{M} \right\}. \quad (23)$$

The challenge we face in using Theorem 11 or the generally tighter bound given in [26], to compute finite blocklength converse bounds is the optimization with respect to the distribution on the set of  $n$ -dimensional input vectors.

The Shannon–Gallager–Berlekamp sphere-packing bound [27] is given by the following result.

*Theorem 12 (Shannon–Gallager–Berlekamp):* Let  $P_{Y|X} : \mathcal{A} \mapsto \mathcal{B}$  be a DMC. Then any  $(n, M, \epsilon)$  code (average probability of error) satisfies

$$\epsilon \geq \exp\{-n(E_{sp}(R - o_1) + o_2)\} \quad (24)$$

where

$$R = \frac{\log M}{n} \quad (25)$$

$$E_{sp}(R) = \sup_{\rho \geq 0} [E_0(\rho) - \rho R] \quad (26)$$

$$E_0(\rho) = \max_{P_X} E_0(\rho, P_X) \quad (27)$$

$$E_0(\rho, P_X) = -\log \sum_{y \in \mathcal{B}} \left[ \sum_{x \in \mathcal{A}} P_X(x) P_{Y|X}(y|x)^{1/(1+\rho)} \right]^{1+\rho} \quad (28)$$

$$= -\log \left( \mathbb{E} \left[ \mathbb{E} \left[ \exp \frac{i(\bar{X}; Y)}{1+\rho} \middle| Y \right] \right]^{1+\rho} \right) \quad (29)$$

$$o_1 = \frac{\log 4}{n} + \frac{|\mathcal{A}| \log n}{n} \quad (30)$$

$$o_2 = \sqrt{\frac{8}{n}} \log \frac{e}{\sqrt{P_{\min}}} + \frac{\log 8}{n} \quad (31)$$

$$P_{\min} = \min\{P_{Y|X}(y|x) : P_{Y|X}(y|x) > 0\} \quad (32)$$

where the maximization in (27) is over all probability distributions on  $\mathcal{A}$ ; and in (29),  $\bar{X}$  and  $Y$  are independent

$$P_{\bar{X}Y}(a, b) = P_X(a) \left( \sum_{x \in \mathcal{A}} P_{Y|X}(b|x) P_X(x) \right). \quad (33)$$

While Theorem 12 is of paramount importance in the analysis of the reliability function for sufficiently high rates, its usefulness in the finite-length regime is more limited because of its slackness and slow speed of convergence of its normalized logarithm. References [8] and [11] have provided tightened versions of the sphere-packing converse bound, which also apply to continuous-output channels.

## E. AWGN Bounds

For the AWGN channel, Shannon [3] gave the following result based on packing spherical cones.

*Theorem 13 (Shannon):* Let

$$Y_i = x_i + Z_i \quad (34)$$

where  $Z_i$  are independent and identically distributed (i.i.d.) standard normal random variables. Assume that each codeword satisfies

$$\sum_{i=1}^n x_i^2 = nP. \quad (35)$$

Define for  $0 \leq \theta \leq \pi/2$

$$q_n(\theta) = Q \left( \sqrt{nP} \right) + \frac{n-1}{\sqrt{\pi}} e^{-nP/2} \int_{\theta}^{\pi/2} (\sin \phi)^{n-2} f_n(\sqrt{nP} \cos \phi) d\phi \quad (36)$$

where

$$f_n(x) = \frac{1}{\Gamma((n+1)/2)} \int_0^{\infty} t^{n-1} e^{-t^2 + \sqrt{2}tx} dt. \quad (37)$$

Then, any  $(n, M, \epsilon)$  code satisfies

$$q_n(\theta(M)) \leq \epsilon \quad (38)$$

with  $\theta(M)$  defined as

$$M\Omega_n(\theta(M)) = \Omega_n(\pi) \quad (39)$$

with

$$\Omega_n(\theta) = \frac{2\pi^{(n-1)/2}}{\Gamma((n-1)/2)} \int_0^\theta (\sin \phi)^{n-2} d\phi \quad (40)$$

which is equal to the area of the unit sphere in  $\mathbb{R}^n$  cut out by a cone with semiangle  $\theta$ . Furthermore, there exists an  $(n, M, \epsilon)$  code with

$$\epsilon \leq q_n(\theta(M)) - \frac{M}{\Omega_n(\pi)} \int_0^{\theta(M)} \Omega_n(\phi) \dot{q}_n(\phi) d\phi \quad (41)$$

$$= \frac{\Gamma(n/2)M}{\sqrt{\pi}\Gamma((n-1)/2)} \int_0^{\theta(M)} q_n(\phi) (\sin \phi)^{n-2} d\phi. \quad (42)$$

Tackled in [4]–[6], [8], and [11], the accurate computation of the bounds in Theorem 13 is challenging.

Applying Theorem 7 to the AWGN channel with  $P_{X^n} = \mathcal{N}(0, P\mathbf{I}_n)$  and optimizing over  $r$  and  $\lambda$ , we obtain the following result (see [16, Theorem 7.4.4]).

*Theorem 14 (Gallager, AWGN):* Consider the AWGN channel with unit noise power and input power  $P$ , with capacity

$$C = \frac{1}{2} \log(1 + P). \quad (43)$$

For blocklength  $n$ , every  $0 \leq R \leq C$  and every  $0 < \alpha \leq 1$ , there exists an  $(\exp(nR), n, \epsilon)$  code (maximal probability of error) with

$$\epsilon \leq \left( \frac{2e^{s(R)n\alpha P}}{\bar{\mu}(\alpha)} \right)^2 e^{-nE_r(R)} \quad (44)$$

where

$$E_r(R) = \frac{P}{4\xi} \left[ (\xi + 1) - (\xi - 1) \sqrt{1 + \frac{4\xi}{P(\xi - 1)}} \right] + \frac{1}{2} \log_e \left\{ \xi - \frac{P(\xi - 1)}{2} \left[ \sqrt{1 + \frac{4\xi}{P(\xi - 1)}} - 1 \right] \right\} \quad (45)$$

for  $R \in [R_c, C]$

$$E_r(R) = 1 - \xi + \frac{P}{2} + \frac{1}{2} \log_e \left( \xi - \frac{P}{2} \right) + \frac{1}{2} \log_e \xi - R \log_e 2, \quad (46)$$

for  $R \in [0, R_c]$

$$\xi = \exp(2 \max\{R, R_c\}) \quad (47)$$

$$R_c = \frac{1}{2} \log \left( \frac{1}{2} + \frac{P}{4} + \frac{1}{2} \sqrt{1 + \frac{P^2}{4}} \right) \quad (48)$$

$$\bar{\mu}(\alpha) = \mathbb{P} \left[ 1 - \alpha \leq \frac{1}{n} \chi_n^2 \leq 1 \right] = \int_{n(1-\alpha)}^n \frac{(y/2)^{n/2-1} e^{-y/2}}{2\Gamma(n/2)} dy \quad (49)$$

$$s(R) = \frac{\rho P}{2(1+\rho)^2 \xi} \quad (50)$$

$$\rho = \frac{P}{2\xi} \left[ 1 + \sqrt{1 + \frac{4\xi}{P(\xi - 1)}} \right] - 1. \quad (51)$$

Other bounds on reliability function have appeared recently, e.g., [12]. However, those bounds provide an improvement only for rates well below capacity.

### F. Normal Approximation

The importance of studying the asymptotics of the function  $M^*(n, \epsilon)$  for given  $\epsilon$  was already made evident by Shannon in [28, Th. 12] which states that, regardless of  $\epsilon \in (0, 1)$

$$\log M^*(n, \epsilon) = nC + o(n) \quad (52)$$

where  $C$  is the channel capacity. Using Theorem 9, Wolfowitz [23] showed (52) for the DMC, and improved the  $o(n)$  term to  $O(\sqrt{n})$  in [24]. Weiss [29] showed that for the BSC with crossover probability  $\delta < \frac{1}{2}$

$$\log_2 M^*(n, \epsilon) \leq n(1 - h(\delta)) - Q^{-1}(\epsilon) \sqrt{n} \sqrt{\delta - \delta^2} \log_2 \frac{1 - \delta}{\delta} + o(\sqrt{n}) \quad (53)$$

where  $Q^{-1}$  denotes the functional inverse of the  $Q$ -function. Crediting M. Pinsker for raising the question, a generalization of (53) was put forward without proof by Dobrushin [30], for symmetric DMCs whose transition matrices are such that the rows are permutation of each other and so are the columns. These results were significantly strengthened and generalized by Strassen [31] who showed that for the DMC

$$\log M^*(n, \epsilon) = nC - Q^{-1}(\epsilon) \sqrt{nV} + O(\log n) \quad (54)$$

where  $V$  denotes the variance of the information density  $i(X; Y)$  under the capacity achieving distribution  $P_X$ ; if such distribution is not unique, then, among those distributions that maximize the average of  $i(X; Y)$  we choose the one that minimizes the variance of  $i(X; Y)$  (if  $\epsilon < 1/2$ ) or that maximizes it (if  $\epsilon \geq 1/2$ ). Strassen's approach in [31] is not amenable to generalization to channels with input constraints (most notably, the AWGN channel). In particular, Theorem 1 is not sufficient to prove the counterpart of (54) to the AWGN channel.

## III. NEW BOUNDS ON RATE

### A. Achievability: Random-Coding

The upper bounds on the average probability of error considered in this paper are based on random coding. The first result gives a general, exact analysis of the error probability of the maximum-likelihood decoder averaged over all codes.

*Theorem 15:* (Random coding average error probability) Denote by  $\epsilon(c_1, \dots, c_M)$  the error probability achieved by the maximum likelihood decoder with codebook  $(c_1, \dots, c_M)$ . Let

$X_1, \dots, X_M$  be independent with marginal distribution  $P_X$ . Then

$$\mathbb{E}[\epsilon(X_1, \dots, X_M)] = 1 - \sum_{\ell=0}^{M-1} \binom{M-1}{\ell} \frac{1}{\ell+1} \mathbb{E}[W^\ell Z^{M-1-\ell}] \quad (55)$$

where

$$W = \mathbb{P}[i(\bar{X}; Y) = i(X; Y) \mid X, Y] \quad (56)$$

$$Z = \mathbb{P}[i(\bar{X}; Y) < i(X; Y) \mid X, Y] \quad (57)$$

with

$$P_{XY\bar{X}}(a, b, c) = P_X(a)P_{Y|X}(b|a)P_X(c). \quad (58)$$

*Proof:* Since the  $M$  messages are equiprobable, upon receipt of the channel output  $y$ , the maximum likelihood decoder chooses with equal probability among the members of the set

$$\arg \max_{i=1, \dots, M} i(c_i; y).$$

Therefore, if the codebook is  $(c_1, \dots, c_M)$ , and  $m = 1$  is transmitted, the maximum likelihood decoder will choose  $\hat{m} = 1$  with probability  $\frac{1}{1+\ell}$  if

$$\sum_{j=2}^M \mathbb{1}\{i(c_j; y) = i(c_1; y)\} = \ell \quad (59)$$

$$\sum_{j=2}^M \mathbb{1}\{i(c_j; y) > i(c_1; y)\} = 0 \quad (60)$$

for  $\ell = 0, \dots, M-1$ . If (60) is not satisfied an error will surely occur. Since the codewords are chosen independently with identical distributions, given that the codeword assigned to message 1 is  $c_1$  and given that the channel output is  $y \in \mathcal{B}$ , the joint distribution of the remaining codewords is  $P_X \times \dots \times P_X$ . Consequently, the conditional probability of correct decision is shown in (61) at the bottom of the page, where  $\bar{X}$  has the same distribution as  $X$ , but is independent of any other random variable arising in this analysis. Averaging (61) with respect to  $(c_1, y)$  jointly distributed as  $P_X P_{Y|X}$  we obtain the summation in (55). Had we conditioned on a message other than  $m = 1$  we would have obtained the same result. Therefore, the error probability averaged over messages and codebook is given by (55). ■

Naturally, Theorem 15 leads to an achievability upper bound since there must exist an

$(M, \mathbb{E}[\epsilon(X_1, \dots, X_M)])$  (average error probability) code.

## B. Achievability: Random-Coding Union Bound

One way to loosen (55) in order to obtain a simpler bound is via the following result.

*Theorem 16: (RCU bound)* For an arbitrary  $P_X$  there exists an  $(M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq \mathbb{E}[\min\{1, (M-1)\mathbb{P}[i(\bar{X}; Y) \geq i(X; Y) \mid X, Y]\}] \quad (62)$$

where  $P_{XY\bar{X}}(a, b, c) = P_X(a)P_{Y|X}(b|a)P_X(c)$ .

*Proof:*<sup>5</sup> The average probability of error attained by an arbitrary codebook  $(c_1, \dots, c_M)$  using a maximum likelihood decoder is upper bounded by

$$\epsilon \leq \frac{1}{M} \sum_{m=1}^M \mathbb{P}\left[\bigcup_{j=1; j \neq m}^M \{i(c_j; Y) \geq i(c_m; Y)\} \mid X = c_m\right] \quad (63)$$

where we do not necessarily have equality since the maximum likelihood decoder resolves some ties in favor of the correct codeword. Using Shannon's random coding argument, the desired result will follow if we can show that the expectation of the RHS of (63) is upper bounded by the RHS of (62) when the codebook is chosen by independent drawings from  $P_X$ . The expectations of all of the  $M$  terms in (63) are identical and are equal to

$$\begin{aligned} & \mathbb{P}\left[\bigcup_{j=2}^M \{i(X_j; Y) \geq i(X_1; Y)\}\right] \\ &= \mathbb{E}\left[\mathbb{P}\left[\bigcup_{j=2}^M \{i(X_j; Y) \geq i(X_1; Y)\} \mid X_1, Y\right]\right] \quad (64) \\ &\leq \mathbb{E}[\min\{1, (M-1)\mathbb{P}[i(X_2; Y) \geq i(X_1; Y) \mid X_1, Y]\}] \quad (65) \end{aligned}$$

where (64) holds by conditioning and averaging, (65) holds by choosing the tighter bound on probability between 1 and the union bound, and all probabilities are with respect to the distribution

$$P_{YX_1 \dots X_M}(b, a_1, \dots, a_M) = P_{Y|X}(b|a_1) \prod_{i=1}^M P_{X_i}(a_i). \quad (66)$$

The proof is now complete since the RHS of (65) is equal to the RHS of (62). ■

<sup>5</sup>A short proof of Theorem 16 can be obtained as a corollary of Theorem 15 by keeping only the first term in the sum (55) and then further upper-bounding  $1 - Z^{M-1}$  by  $\min\{(M-1)(1-Z), 1\}$ . The standalone proof we give here is useful in Appendix A.

$$\sum_{\ell=0}^{M-1} \binom{M-1}{\ell} \frac{1}{\ell+1} (\mathbb{P}[i(\bar{X}; y) = i(c_1; y)])^\ell (\mathbb{P}[i(\bar{X}; y) < i(c_1; y)])^{M-1-\ell} \quad (61)$$

Gallager’s bound (7) can also be obtained by analyzing the average behavior of random coding and maximum-likelihood decoding. In fact, it is easy to verify that we can weaken (62) to recover (7) using  $\max\{0, x\} \leq x^{1/(1+\lambda)}$  and  $\min\{x, 1\} \leq x^\lambda$ . Furthermore Shannon’s bound (6) can also be obtained by weakening (62) by splitting the expectation according to whether or not  $i(X; Y) \leq \log \beta$  and upper bounding  $\min\{x, 1\}$  by 1 when  $i(X; Y) \leq \log \beta$  and by  $x$  otherwise.

In principle, without exploiting any symmetries, the brute-force computation of the bound (62) has complexity  $O(n^{2(|\mathcal{A}|-1)|\mathcal{B}|})$  for a DMC with input/output alphabets  $\mathcal{A}$  and  $\mathcal{B}$ . Next we give easier-to-compute upper bounds that do not sacrifice much tightness.

*C. Achievability: Dependence Testing Bound*

*Theorem 17:* (DT bound) For any distribution  $P_X$  on  $\mathcal{A}$ , there exists a code with  $M$  codewords and average probability of error not exceeding

$$\epsilon \leq \mathbb{E} \left[ \exp \left\{ - \left[ i(X; Y) - \log \frac{M-1}{2} \right]^+ \right\} \right]. \quad (67)$$

*Proof:* Consider the following obvious identity for arbitrary  $z \geq 0$  and  $\gamma > 0$ :

$$\exp \left\{ - \left[ \log \frac{z}{\gamma} \right]^+ \right\} = 1\{z \leq \gamma\} + \frac{\gamma}{z} 1\{z > \gamma\} \quad (68)$$

(for  $z = 0$  we understand both sides to be equal to 1, which is the value attained for all  $0 < z \leq \gamma$ ). If we let  $z = \frac{dP_{XY}}{d(P_X \times P_Y)}$  and we average both sides of (68) with respect to  $P_{XY}$  we obtain

$$\mathbb{E} \left[ \exp \left\{ - [i(X; Y) - \log \gamma]^+ \right\} \right] = \mathbb{P}[i(X; Y) \leq \log \gamma] + \gamma \mathbb{P}[i(X; \bar{Y}) > \log \gamma]. \quad (69)$$

Letting  $\gamma = \frac{M-1}{2}$ , we see that Theorem 17 is, in fact, equivalent to the following result. ■

*Theorem 18:* For any distribution  $P_X$  on  $\mathcal{A}$ , there exists a code with  $M$  codewords and average probability of error not exceeding

$$\epsilon \leq \mathbb{P} \left[ i(X; Y) \leq \log \frac{M-1}{2} \right] + \frac{M-1}{2} \mathbb{P} \left[ i(X; \bar{Y}) > \log \frac{M-1}{2} \right] \quad (70)$$

where  $P_{XY\bar{Y}}(a, b, c) = P_X(a)P_{Y|X}(b|a)P_Y(c)$ .

*Proof:* The proof combines Shannon’s random coding with Feinstein’s suboptimal decoder. Fix  $P_X$ . Let  $\{Z_x\}_{x \in \mathcal{A}} : \mathcal{B} \mapsto \{0, 1\}$  be a collection of deterministic functions defined as

$$Z_x(y) = 1 \left\{ i(x; y) > \log \frac{M-1}{2} \right\}. \quad (71)$$

For a given codebook  $(c_1, \dots, c_M)$ , the decoder runs  $M$  likelihood ratio binary hypothesis tests in parallel, the  $j^{\text{th}}$  of which is between the true distribution  $P_{Y|X=c_j}$  and “average noise”  $P_Y$ . The decoder computes the values  $Z_{c_j}(y)$  for the received

channel output  $y$  and returns the lowest index  $j$  for which  $Z_{c_j}(y) = 1$  (or declares an error if there is no such index). The conditional error probability given that the  $j$ th message was sent is

$$\begin{aligned} & \mathbb{P} \left[ \{Z_{c_j}(Y) = 0\} \cup_{i < j} \{Z_{c_i}(Y) = 1\} \mid X = c_j \right] \\ & \leq \mathbb{P}[i(c_j; Y) \leq \log \frac{M-1}{2} \mid X = c_j] \\ & \quad + \sum_{i < j} \mathbb{P}[i(c_i; Y) > \log \frac{M-1}{2} \mid X = c_j] \end{aligned} \quad (72)$$

where we have used the union bound and the definition of  $Z_x(y)$ . Averaging (72) over codebooks  $\{c_i\}$  that are generated as (pairwise) independent random variables with distribution  $P_X$  we obtain

$$\mathbb{P} \left[ i(X; Y) \leq \log \frac{M-1}{2} \right] + (j-1) \mathbb{P} \left[ i(X; \bar{Y}) > \log \frac{M-1}{2} \right]$$

where recall that  $\bar{Y}$  has the same distribution as  $Y$ , but unlike  $Y$ , it is independent of  $X$ . Averaging further over equiprobable messages, and since

$$\frac{1}{M} \sum_{j=1}^M (j-1) = \frac{M-1}{2} \quad (73)$$

we obtain that the average error probability is upper bounded by (70), and therefore there must exist a code whose average error probability is upper bounded by that expression. ■

We may wonder whether in the above proof a choice of threshold different from  $\log \frac{M-1}{2}$  may lead to a tighter bound. In fact, it is readily seen that we can generalize Theorem 18 not just to any other constant value of the threshold but to thresholds that are codeword dependent, leading to the following result.

*Lemma 19:* For any distribution  $P_X$  on  $\mathcal{A}$ , and any measurable function  $\gamma : \mathcal{A} \mapsto [0, \infty]$ , there exists an  $(M, \epsilon)$  code (average probability of error) satisfying

$$\epsilon \leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] + \frac{M-1}{2} \mathbb{P}[i(X; \bar{Y}) > \log \gamma(X)] \quad (74)$$

where  $P_{XY\bar{Y}}(a, b, c) = P_X(a)P_{Y|X}(b|a)P_Y(c)$ .

In order to optimize the choice of the function  $\gamma(\cdot)$ , we can view the RHS of (74) as the average with respect to  $P_X$  of

$$P_{Y|X=x}[i(x; Y) \leq \log \gamma(x)] + \frac{M-1}{2} P_Y[i(x; Y) > \log \gamma(x)] \quad (75)$$

which is a weighted sum of two types of errors. Thus, for every  $x$ , (75) is equal to  $\frac{M+1}{2}$  times the average error probability in a Bayesian hypothesis testing problem between  $P_{Y|X=x}$  with *a priori* probability  $\frac{2}{M+1}$  and  $P_Y$  with *a priori* probability  $\frac{M-1}{M+1}$ . The average error probability is then minimized by the test that compares the likelihood ratio between these two distributions to the ratio of the two *a priori* probabilities. Thus, we obtain that the optimal threshold is, in fact, codeword independent:  $\gamma(x) = \frac{M-1}{2}$ ; and Theorem 18 gives the tightest version of Lemma 19.

*Remarks:*

- 1) Unlike the existing bounds (5), (6), and (7), the bounds in Theorems 17 and 18 require no selection of auxiliary constants.
- 2) Theorem 2 follows by taking  $\gamma(x) = \beta$  in Lemma 19 and weakening  $\frac{M-1}{2}$  by a factor of 2.
- 3) The bound in [32] is provably weaker than Theorem 17 (originally published in [33]).
- 4) It can be easily seen from (67) that Theorem 17 can be used to prove the achievability part of the most general known channel capacity formula [2].
- 5) We refer to the bound in Theorems 17 and 18 as the *dependence testing* bound because the RHS of (70) is equal to  $\frac{M+1}{2}$  times the Bayesian minimal error probability of a binary hypothesis test of dependence:

$$\begin{aligned} H_1 : P_{XY} & \text{ with probability } \frac{2}{M+1} \\ H_0 : P_X P_Y & \text{ with probability } \frac{M-1}{M+1} \end{aligned}$$

- 6) An alternative expression for (67) is given by

$$\epsilon \leq \int_0^1 \mathbb{P} \left[ i(X; Y) \leq \log \frac{M-1}{2u} \right] du. \quad (76)$$

This follows from (67) and:

$$\mathbb{E} [\min\{X, 1\}] = \int_0^1 \mathbb{P}[X \geq u] du \quad (77)$$

which is valid for any nonnegative  $X$ .

- 7) Yet another way to look at (67) is by defining a particular  $f$ -divergence [34] as follows:

$$\mathcal{D}_M(P||Q) = \int \left[ \frac{dP}{dQ} - \frac{M-1}{2} \right]^+ dQ. \quad (78)$$

Then (67) is equivalent to

$$1 - \epsilon \geq \mathcal{D}_M(P_{XY}||P_X P_Y). \quad (79)$$

Since processing does not increase  $f$ -divergence, the lower bound (79) can be further simplified by applying a suitable mapping of the space  $A \times B$  into some other space.

Using Lemma 19 we can easily extend Theorem 18 to the case of input constraints.

*Theorem 20:* For any distribution  $P_X$  on  $A$  there exists a code with  $M$  codewords in  $F$  and average probability of error satisfying

$$\begin{aligned} \epsilon \leq \mathbb{P} \left[ i(X; Y) \leq \log \frac{M-1}{2} \right] \\ + \frac{M-1}{2} \mathbb{P} \left[ i(X; \bar{Y}) > \log \frac{M-1}{2} \right] + P_X[F^c]. \end{aligned} \quad (80)$$

*Proof:* Set  $\gamma(x) = \frac{M-1}{2}$  for  $x \in F$  and  $\gamma(x) = +\infty$  for  $x \in F^c$ . Then by Lemma 19 we have

$$\begin{aligned} \epsilon \leq \mathbb{P} \left[ \left\{ i(X; Y) \leq \log \frac{M-1}{2} \right\} \cup \left\{ X \in F^c \right\} \right] \\ + \frac{M-1}{2} \mathbb{P} \left[ i(X; \bar{Y}) > \log \frac{M-1}{2}, X \in F \right]. \end{aligned} \quad (81)$$

Trivial upper bounding yields (80). Lemma 19 guarantees the existence of a codebook whose average probability of error satisfies the required (80). However, we are not guaranteed that that codebook is feasible since some of the codewords might fall outside the set  $F$ . If we modify the codebook, replacing every infeasible codeword by an arbitrary  $c_0 \in F$ , while not modifying the decoder, the error probability (averaged over messages) does not change. The reason is that the decoding set corresponding to a message that has been assigned an infeasible codeword is empty (because the corresponding threshold is  $+\infty$ ), and therefore, its conditional probability of error is 1, and remains 1 after it has been replaced by  $c_0$  since the decoder has not been modified. ■

#### D. Achievability: Maximal Probability of Error

Any achievability bound on average error probability gives a bound on maximal error probability since the existence of an  $(M, \epsilon)$  code in the average sense guarantees the existence of an  $(M - \frac{M\epsilon}{\epsilon'}, \epsilon')$  code in the maximal sense, for any  $\epsilon < \epsilon' < 1$ . However, in this subsection we give maximal error probability counterparts to some of the bounds in Section III-C.

1) *Bounds Fixing the Input Distribution:* As we saw in the proof of Theorem 18, the random coding method is such that only pairwise independent codewords are required. If  $A = \mathcal{A}^n$ ,  $M = |\mathcal{A}|^k$ , and  $\mathcal{A}$  is a finite field, then an interesting ensemble that satisfies that property (but not total statistical independence) together with  $P_X$  being equiprobable on  $A$  is that of a random linear code: construct a random  $n \times k$  matrix with independent coefficients equiprobable on  $\mathcal{A}$ ; then the  $M$  codewords are generated as the products of the matrix and every vector in  $\mathcal{A}^k$ . For certain channels such as additive-noise discrete channels and erasure channels, the average error probability and the maximal error probability coincide for linear codes (with an appropriately defined randomized maximum likelihood (ML) decoder; see Appendix A). Therefore, for those channels, the bound in Theorem 17 achieved with an equiprobable  $P_X$  not only can be achieved by a linear code but it is also an upper bound on maximal error probability.

The following bound on maximal error probability holds in general.

*Theorem 21:* For any input distribution  $P_X$  and measurable  $\gamma : A \rightarrow [0, \infty]$ , there exists a code with  $M$  codewords such that the  $j$ th codeword's probability of error satisfies

$$\begin{aligned} \epsilon_j \leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] \\ + (j-1) \sup_x \mathbb{P}[i(x; Y) > \log \gamma(x)] \end{aligned} \quad (82)$$

where the first probability is with respect to  $P_{XY}$  and the second is with respect to the unconditional distribution  $P_Y$ . In particular, the maximal probability of error satisfies

$$\begin{aligned} \epsilon \leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] \\ + (M-1) \sup_x \mathbb{P}[i(x; Y) > \log \gamma(x)]. \end{aligned} \quad (83)$$

*Proof:* First, we specify the operation of the decoder given the codebook  $(c_1, \dots, c_M)$ . The decoder simply computes



$i(c_j; y)$  for the received channel output  $y$  and selects the first codeword  $c_j$  for which  $i(c_j; y) > \log \gamma(c_j)$ .

Now, let us show that we can indeed choose  $M$  codewords so that their respective probabilities of decoding error satisfy (82). For the first codeword, the conditional probability of error under the specified decoding rule is independent of other codewords and is equal to

$$\epsilon_1(x) = \mathbb{P}[i(x; Y) \leq \log \gamma(x) | X = x] \quad (84)$$

if the first codeword is  $x \in A$ . There must exist at least one choice of  $x$ , which we call  $c_1$ , such that

$$\epsilon_1(c_1) \leq \mathbb{E}[\epsilon_1(X)] \quad (85)$$

$$= \mathbb{P}[i(X; Y) \leq \log \gamma(X)]. \quad (86)$$

Now assume that  $j - 1$  codewords  $\{c_\ell\}_{\ell=1}^{j-1}$  have been chosen and we are to show that  $c_j$  can also be chosen so that (82) is satisfied. Denote

$$D_{j-1} = \bigcup_{\ell=1}^{j-1} \{y : i(c_\ell; y) > \log \gamma(c_\ell)\} \subseteq B. \quad (87)$$

If the  $j$ th codeword is  $x$ , its conditional probability of error is

$$\begin{aligned} & \epsilon_j(c_1, \dots, c_{j-1}, x) \\ &= 1 - \mathbb{P}[\{i(x; Y) > \log \gamma(x)\} \setminus D_{j-1} | X = x]. \end{aligned} \quad (88)$$

Thus

$$\begin{aligned} & \mathbb{E}[\epsilon_j(c_1, \dots, c_{j-1}, X)] \\ &= \mathbb{P}[\{i(X; Y) \leq \log \gamma(X)\} \cup D] \end{aligned} \quad (89)$$

$$\leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] + P_Y(D) \quad (90)$$

$$\leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] + (j - 1) \sup_{x \in A} P_Y[i(x; Y) > \log \gamma(x)]. \quad (91)$$

Thus, there must exist a codeword  $c_j$  such that  $\epsilon_j(c_1, \dots, c_{j-1}, c_j)$  satisfies (82). ■

By upper-bounding the second term in (83) via

$$\mathbb{P}[i(x; Y) \geq \log \gamma] \leq \frac{1}{\gamma} \quad (92)$$

we observe that Feinstein's Theorem 1 is a corollary of Theorem 21.

The proof technique we used to show Theorem 21 might be called sequential random coding because each codeword is chosen sequentially depending on the previous choices, and its existence is guaranteed by the fact that the average cannot be exceeded by every realization. Note that there is no contradiction due to the nonending nature of sequential random coding: sooner or later the conditional probability of error of the next message becomes 1.

Some symmetric channels and choices of  $P_X$  (most notably the BEC and the BSC under equiprobable  $P_X$ ) satisfy the sufficient condition in the next result.

*Theorem 22:* Fix an arbitrary input distribution  $P_X$ . If the cumulative distribution function  $\mathbb{P}[i(x; Y) \leq \alpha]$  does not depend

on  $x$  for any  $\alpha$  when  $Y$  is distributed according to  $P_Y$ , then there exists an  $(M, \epsilon)$  code with maximal probability of error satisfying

$$\epsilon \leq \mathbb{E} \left[ \exp \left\{ - [i(X; Y) - \log(M - 1)]^+ \right\} \right]. \quad (93)$$

*Proof:* Under the stated conditions, (83) states that the maximal error probability is upper bounded by the average with respect to  $x \sim P_X$  of

$$P_{Y|X=x}[i(x; Y) \leq \log \gamma(x)] + (M - 1) P_Y[i(x; Y) > \log \gamma(x)]. \quad (94)$$

Thus,  $\gamma(x)$  can be optimized similarly to (75). ■

2) *Extension to Input Constraints:* Theorem 21 can be extended to the case of input constraints in the following way.

*Theorem 23:* For any input distribution  $P_X$  and measurable  $\gamma : A \rightarrow [0, \infty]$ , there exists a code with  $M$  codewords in the set  $F$  such that the maximal probability of error  $\epsilon$  satisfies

$$\begin{aligned} \epsilon P_X[F] &\leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] \\ &+ (M - 1) \sup_{x \in F} \mathbb{P}[i(x; Y) > \log \gamma(x)]. \end{aligned} \quad (95)$$

*Proof:* The proof is the same as that of Theorem 21 with the modification that the selection of each codeword belongs to  $F$ , and at each step we use the fact that for an arbitrary nonnegative function  $g : A \mapsto \mathbb{R}$ , there exists  $x \in F$  such that

$$g(x) \leq \frac{\mathbb{E}[g(X)]}{P_X[F]} \quad (96)$$

since otherwise we would get the impossible  $\mathbb{E}[g(X)] < \mathbb{E}[1\{X \in F\}g(X)]$ . ■

Comparing Theorem 23 with Theorem 20 we note that (95) is stronger than the bound

$$\begin{aligned} \epsilon &\leq \mathbb{P}[i(X; Y) \leq \log \gamma(X)] \\ &+ (M - 1) \sup_{x \in F} \mathbb{P}[i(x; Y) > \log \gamma(x)] + P_X[F^c]. \end{aligned} \quad (97)$$

Using the fact that

$$P_Y[i(x; Y) \geq \log \gamma] \leq \frac{1}{\gamma} \quad (98)$$

an immediate corollary of Theorem 23 is the following.

*Theorem 24:* For any distribution  $P_X$  and any  $\gamma > 0$ , there exists an  $(M, \epsilon)$  code (maximal probability of error) with codewords in the set  $F \subset A$  such that

$$M \geq 1 + \gamma (\epsilon P_X[F] - \mathbb{P}[i(X; Y) < \log \gamma]). \quad (99)$$

Note that (99) is always stronger than the conventional input-constrained version of Feinstein's bound (14).

3) *Bounds Fixing the Output Distribution:* All the previous achievability bounds fixed some input distribution  $P_X$  and then proved that a certain codebook exists. However, in some cases (most notably, the AWGN channel) it is desirable to consider

auxiliary distributions on the output alphabet that are not necessarily induced by an input distribution.

The optimal performance of binary hypothesis testing plays an important role in our development. Consider a random variable  $W$  defined on  $W$  which can take probability measures  $P$  or  $Q$ . A randomized test between those two distributions is defined by a random transformation  $P_{Z|W} : W \mapsto \{0, 1\}$  where 0 indicates that the test chooses  $Q$ . The best performance achievable among those randomized tests is given by<sup>6</sup>

$$\beta_\alpha(P, Q) = \min_{P_{Z|W}: \sum_{w \in W} P(w)P_{Z|W}(1|w) \geq \alpha} \sum_{w \in W} Q(w)P_{Z|W}(1|w) \quad (100)$$

where the minimum is guaranteed to be achieved by the Neyman–Pearson lemma (Appendix B). Thus,  $\beta_\alpha(P, Q)$  gives the minimum probability of error under hypothesis  $Q$  if the probability of error under hypothesis  $P$  is not larger than  $1 - \alpha$ . As a function of  $\alpha$ , (100) is a piecewise-linear convex function joining the points

$$\begin{aligned} \beta_\alpha &= \mathbb{Q} \left[ \frac{dP}{dQ} \geq \gamma \right] \\ \alpha &= \mathbb{P} \left[ \frac{dP}{dQ} \geq \gamma \right] \end{aligned} \quad (101)$$

iterated over all  $\gamma > 0$ . It is easy to show that (e.g., [35]) for any  $\gamma > 0$

$$\alpha \leq \mathbb{P} \left[ \frac{dP}{dQ} \geq \gamma \right] + \gamma \beta_\alpha(P, Q). \quad (102)$$

On the other hand

$$\beta_\alpha(P, Q) \leq \frac{1}{\gamma_0} \quad (103)$$

where  $\gamma_0$  satisfies

$$\mathbb{P} \left[ \frac{dP}{dQ} \geq \gamma_0 \right] \geq \alpha. \quad (104)$$

Additional results on the behavior of  $\beta_\alpha(P, Q)$  in the case when  $P$  and  $Q$  are product distributions are given in Appendix C.

Throughout most of our development, the binary hypothesis testing of interest is  $W = B$ ,  $P = P_{Y|X=x}$  and  $Q = Q_Y$ , an auxiliary unconditional distribution.<sup>7</sup> In that case, for brevity and with a slight abuse of notation we will denote

$$\beta_\alpha(x, Q_Y) = \beta_\alpha(P_{Y|X=x}, Q_Y). \quad (105)$$

As a consequence of (102) we have

$$\beta_\alpha(x, Q_Y) \geq \sup_{\gamma > 0} \frac{1}{\gamma} \left( \alpha - P_{Y|X=x} \left[ \frac{dP_{Y|X=x}}{dQ_Y} \geq \gamma \right] \right). \quad (106)$$

<sup>6</sup>We write summations over alphabets for simplicity; however, all of our general results hold for arbitrary probability spaces.

<sup>7</sup>As we show later, it is sometimes advantageous to allow  $Q_Y$  that cannot be generated by any input distribution.

Each per-codeword cost constraint can be defined by specifying a subset  $F \subset A$  of permissible inputs. For an arbitrary  $F \subset A$ , we define a related measure of performance for the composite hypothesis test between  $Q_Y$  and the collection  $\{P_{Y|X=x}\}_{x \in F}$

$$\kappa_\tau(F, Q_Y) = \inf_{P_{Z|Y}: \inf_{x \in F} P_{Z|X}(1|x) \geq \tau} \sum_{y \in B} Q_Y(y)P_{Z|Y}(1|y). \quad (107)$$

Again, typically we will take  $A$  and  $B$  as  $n$ -fold Cartesian products of alphabets  $\mathcal{A}$  and  $\mathcal{B}$ . To emphasize dependence on  $n$  we will write  $\beta_\alpha^n(x, Q_Y)$  and  $\kappa_\tau^n(F, Q_Y)$ . Since  $Q_Y$  and  $F$  will usually be fixed we will simply write  $\kappa_\tau^n$ . Also, in many cases  $\beta_\alpha^n(x, Q_Y)$  will be the same for all  $x \in F$ . In these cases we will write  $\beta_\alpha^n$ .

*Theorem 25 (Achievability, Input Constraints:  $\kappa\beta$  Bound):* For any  $0 < \epsilon < 1$ , there exists an  $(M, \epsilon)$  code with codewords chosen from  $F \subset A$ , satisfying

$$M \geq \sup_{0 < \tau < \epsilon} \sup_{Q_Y} \frac{\kappa_\tau(F, Q_Y)}{\sup_{x \in F} \beta_{1-\epsilon+\tau}(x, Q_Y)}. \quad (108)$$

*Note:* It is possible<sup>8</sup> that (108) will be of the form  $M \geq \alpha/0$  with  $\alpha > 0$ . In this case the statement of the theorem should be understood as “ $(M, \epsilon)$  codes with arbitrarily high  $M$  exist.”

*Proof:* Fix  $0 < \epsilon < 1$ ,  $0 < \tau < \epsilon$ , and  $Q_Y$ . We construct the collection of random binary-valued  $\{Z_x\}_{x \in F}$  conditionally independent given  $Y$ , and with marginal conditional distributions given by  $P_{Z_x|Y}$ , which denotes, for brevity, the conditional distribution that achieves the minimum in (100) for  $\alpha = 1 - \epsilon + \tau$ ,  $Q_Y$ , and  $x \in F$ .

We construct the codebook sequentially:

Step 1. Choose  $c_1 \in F$  arbitrarily. Note that regardless of the choice of  $c_1$  we have from (100) that

$$\mathbb{P}[Z_{c_1} = 1|X = c_1] \geq 1 - \epsilon + \tau. \quad (109)$$

Step  $k$ . Assume  $c_1, \dots, c_{k-1}$  have been chosen. Choose  $c_k \in F$  so that

$$\begin{aligned} &\mathbb{P}[Z_{c_1} = 0, \dots, Z_{c_{k-1}} = 0, Z_{c_k} = 1|X = c_k] \\ &= \sum_{y \in B} P_{Y|X}(y|c_k)P_{Z_{c_k}|Y}(1|y) \prod_{i=1}^{k-1} P_{Z_{c_i}|Y}(0|y) \\ &> 1 - \epsilon. \end{aligned} \quad (110)$$

Unless such a choice is impossible, proceed to the next step.

Let  $M$  be the number of steps that this procedure takes before stopping. (In case it does not stop, we let  $M = \infty$ .)

The decoder simply applies the  $M$  independent random transformations  $P_{Z_{c_1}|Y}^*, \dots, P_{Z_{c_M}|Y}^*$  to the data. If all  $M$  outputs are 0, the decoder outputs 0; otherwise, it outputs the smallest  $i$  such that  $Z_{c_i} = 1$ .

<sup>8</sup>For an example of such a case, take  $A = B = [0, 1]$  with the Borel  $\sigma$ -algebra. Define  $P_{Y|X=x}(y) = \delta_x(y)$ , i.e., a point measure at  $y = x$ , and take  $Q_Y$  to be Lebesgue measure. Then,  $\beta_\alpha(x, Q_Y) = 0$  for any  $x$  and  $\alpha$ , and  $\kappa_\tau(Q_Y) = 1$  for any  $\tau > 0$ .

It follows from the encoder/decoder construction and (110) that the maximal error probability of the code is indeed upper bounded by  $\epsilon$ . Let

$$Z^* = Z_{c_1} \text{ OR } \dots \text{ OR } Z_{c_M}. \quad (111)$$

For any  $x \in \mathbb{F}$ , we have

$$P_{Z^*|X}(0|x) = \mathbb{P}[Z^* = 0, Z_x = 1|X = x] + \mathbb{P}[Z^* = 0, Z_x = 0|X = x] \quad (112)$$

$$\leq 1 - \epsilon + \mathbb{P}[Z^* = 0, Z_x = 0|X = x] \quad (113)$$

$$\leq 1 - \epsilon + \mathbb{P}[Z_x = 0|X = x] \quad (114)$$

$$< 1 - \epsilon + \epsilon - \tau \quad (115)$$

where (113) follows because if  $x \in \{c_1, \dots, c_M\}$ , it is impossible that  $Z^* = 0$  and  $Z_x = 1$  simultaneously, while if  $x \in \mathbb{F} - \{c_1, \dots, c_M\}$  we were not able to add  $x$  to the codebook, and therefore  $\mathbb{P}[Z^* = 0, Z_x = 1|X = x] \leq 1 - \epsilon$ ; and (115) follows by the construction of  $Z_x$  from (100).

From (115) we conclude that  $P_{Z^*|X}$  is such that

$$\inf_{x \in \mathbb{F}} P_{Z^*|X}(1|x) \geq \tau. \quad (116)$$

Accordingly

$$\kappa_\tau(\mathbb{F}, Q_Y) \leq \sum_{y \in \mathbb{B}} Q_Y(y) P_{Z^*|Y}(1|y) \quad (117)$$

$$\leq \sum_{y \in \mathbb{B}} Q_Y(y) \sum_{m=1}^M P_{Z_{c_m}|Y}(1|y) \quad (118)$$

$$= \sum_{m=1}^M \beta_{1-\epsilon+\tau}(c_m, Q_Y) \quad (119)$$

$$\leq M \sup_{x \in \mathbb{F}} \beta_{1-\epsilon+\tau}(x, Q_Y) \quad (120)$$

where (117) follows from (116) and (107); (118) follows from (111); and (119) follows from the fact that by definition of  $P_{Z_{c_m}|Y}$ , it achieves the minimum in (100) for  $x = c_m$  and  $\alpha = 1 - \epsilon + \tau$ . ■

In (100) and (107) we have defined  $\beta_\alpha$  and  $\kappa_\tau$  using randomized tests. Then, in Theorem 25 we have constructed the coding scheme with a randomized decoder. Correspondingly, if we define  $\beta_\alpha$  and  $\kappa_\tau$  using nonrandomized tests, then the analog of Theorem 25 for a nonrandomized decoder can be proved.

As long as  $Q_Y$  is the output distribution induced by an input distribution  $Q_X$ , the quantity (107) satisfies the bounds

$$\tau Q_X[\mathbb{F}] \leq \kappa_\tau(\mathbb{F}, Q_Y) \quad (121)$$

$$\leq \tau. \quad (122)$$

The bound (122) is achieved by choosing the test  $Z$  that is equal to 1 with probability  $\tau$  regardless of  $Y$ ; since  $\kappa_\tau$  is achieved by

the optimal test, it can only be better. To verify (121), note that for any  $P_{Z|Y}$  that satisfies the condition in (107), we have

$$\sum_{y \in \mathbb{B}} Q_Y(y) P_{Z|Y}(1|y) = \sum_{x \in \mathbb{A}} \sum_{y \in \mathbb{B}} Q_X(x) P_{Y|X}(y|x) P_{Z|Y}(1|y) \quad (123)$$

$$\geq \sum_{x \in \mathbb{F}} Q_X(x) \sum_{y \in \mathbb{B}} P_{Y|X}(y|x) P_{Z|Y}(1|y) \quad (124)$$

$$\geq \sum_{x \in \mathbb{F}} Q_X(x) \left\{ \inf_{x \in \mathbb{F}} \sum_{y \in \mathbb{B}} P_{Y|X}(y|x) P_{Z|Y}(1|y) \right\} \quad (125)$$

$$\geq \tau Q_X[\mathbb{F}]. \quad (126)$$

Using (121) in Theorem 25 we obtain a weakened but useful bound:

$$M \geq \sup_{0 < \tau < \epsilon} \sup_{Q_X} \frac{\tau Q_X[\mathbb{F}]}{\sup_{x \in \mathbb{F}} \beta_{1-\epsilon+\tau}(x, Q_Y)} \quad (127)$$

where the supremum is over all input distributions, and  $Q_Y$  denotes the distribution induced by  $Q_X$  on the output. By a judicious choice of  $\gamma(x)$  in Lemma 19 we can obtain a strengthened version of the bound for average error probability with the supremum in the denominator of (127) replaced by the average.

### E. General Converse: Average Probability of Error

We give first a general result, which upon particularization leads to a new converse as well to the recovery of previously known converses; see Section III-G. The statement of the result uses the notation introduced in (100) particularized to  $\mathbb{W} = \mathbb{A} \times \mathbb{B}$ .

*Theorem 26:* For a given code (possibly randomized encoder and decoder pair), let

$\epsilon$  = average error probability with  $P_{Y|X}$

$\epsilon'$  = average error probability with  $Q_{Y|X}$  and

$P_X = Q_X$  = encoder output distribution

with equiprobable codewords.

Then

$$\beta_{1-\epsilon}(P_{XY}, Q_{XY}) \leq 1 - \epsilon'. \quad (128)$$

*Proof:* The message is denoted by the random variable  $S$ , equiprobable on  $\{1, \dots, M\}$ . The encoder and decoder are the random transformations  $P_{X|S}$  and  $P_{Z|Y}$ . Consider the following (suboptimal) test for deciding between  $P_{XY}$  and  $Q_{XY}$ : denote the observed pair by  $(x, y)$ ;  $y$  is fed to the decoder which selects  $z$ , and the test declares  $P_{XY}$  with probability

$$P_{S|X}(z|x) = \frac{P_{X|S}(x|z)}{M P_X(x)}. \quad (129)$$

The probability that the test is correct if  $P_{XY}$  is the actual distribution is

$$\begin{aligned} & \sum_{a \in A} \sum_{b \in B} \sum_{z=1}^M P_X(a) P_{Y|X}(b|a) P_{Z|Y}(z|b) \frac{P_{X|S}(a|z)}{M P_X(a)} \\ &= \sum_{z=1}^M \sum_{a \in A} \sum_{b \in B} \frac{1}{M} P_{X|S}(a|z) P_{Y|X}(b|a) P_{Z|Y}(z|b) \quad (130) \\ &= 1 - \epsilon \quad (131) \end{aligned}$$

where (131) is simply the definition of  $\epsilon$ . Likewise, the probability that the test is incorrect if  $Q_{XY}$  is the actual distribution is

$$\begin{aligned} & \sum_{a \in A} \sum_{b \in B} \sum_{z=1}^M P_X(a) Q_{Y|X}(b|a) P_{Z|Y}(z|b) \frac{P_{X|S}(a|z)}{M P_X(a)} \\ &= \sum_{z=1}^M \sum_{a \in A} \sum_{b \in B} \frac{1}{M} P_{X|S}(a|z) Q_{Y|X}(b|a) P_{Z|Y}(z|b) \quad (132) \\ &= 1 - \epsilon' \quad (133) \end{aligned}$$

where (133) is simply the definition of  $\epsilon'$ .

The optimal test that attains the minimum in (100) among all tests such that the probability of corrected decision under  $P_{XY}$  is not less than  $1 - \epsilon$  has a probability of incorrect decision under  $Q_{XY}$  that cannot be larger than (133). ■

Theorem 26 allows one to use any converse for channel  $Q_{Y|X}$  to prove a converse for channel  $P_{Y|X}$ . It has many interesting generalizations (for example, to list-decoding and channels with feedback) and applications, whose study is outside the scope of this paper.

A simple application of Theorem 26 yields the following result.

*Theorem 27 (Converse):* Every  $(M, \epsilon)$  code (average probability of error) with codewords belonging to  $F$  satisfies

$$M \leq \sup_{P_X} \inf_{Q_Y} \frac{1}{\beta_{1-\epsilon}(P_{XY}, P_X \times Q_Y)} \quad (134)$$

where  $P_X$  ranges over all distributions on  $F$ , and  $Q_Y$  ranges over all distributions on  $B$ .

*Proof:* Denote the distribution of the encoder output by  $\bar{P}_X$  and particularize Theorem 26 by choosing  $Q_{Y|X} = Q_Y$  for an arbitrary  $Q_Y$ , in which case we obtain  $\epsilon' = 1 - \frac{1}{M}$ . Therefore, from (128) we obtain

$$\frac{1}{M} \geq \sup_{Q_Y} \beta_{1-\epsilon}(\bar{P}_X P_{Y|X}, \bar{P}_X \times Q_Y) \quad (135)$$

$$\geq \inf_{P_X} \sup_{Q_Y} \beta_{1-\epsilon}(P_{XY}, P_X \times Q_Y). \quad (136)$$

As we will see shortly in important special cases,  $\beta_\alpha(x, Q_Y)$  is constant on  $F$ . In those cases the following converse is particularly useful.

*Theorem 28:* Fix a probability measure  $Q_Y$  on  $B$ . Suppose that  $\beta_\alpha(x, Q_Y) = \beta_\alpha(Q_Y)$  for  $x \in F$ . Then every  $(M, \epsilon)$ -code (average probability of error) satisfies

$$M \leq \frac{1}{\beta_{1-\epsilon}(Q_Y)}. \quad (137)$$

*Proof:* The result follows from Theorem 27 and the following auxiliary result. ■

*Lemma 29:* Suppose that  $\beta_\alpha(P_{Y|X=x}, Q_{Y|X=x}) = \beta_\alpha$  is independent of  $x \in F$ . Then, for any  $P_X$  supported on  $F$  we have

$$\beta_\alpha(P_X P_{Y|X}, P_X Q_{Y|X}) = \beta_\alpha(P_{Y|X=x}, Q_{Y|X=x}). \quad (138)$$

*Proof:* Take a collection of optimal tests  $Z_x$  for each pair  $P_{Y|X=x}$  versus  $Q_{Y|X=x}$ , i.e.

$$P_{Y|X=x}[Z_x = 1] \geq 1 - \alpha \quad (139)$$

$$Q_{Y|X=x}[Z_x = 1] = \beta_{1-\alpha}. \quad (140)$$

Then take  $Z_X$  as a test for  $P_{XY}$  versus  $Q_{XY}$ . In this way, we get

$$\beta_\alpha(P_X P_{Y|X}, P_X Q_{Y|X}) \leq \beta_\alpha(P_{Y|X=x}, Q_{Y|X=x}). \quad (141)$$

We now prove the reverse inequality. Consider an arbitrary test  $Z$  such that

$$P_{XY}[Z = 1] = \sum_{x \in A} P_X(x) P_{Y|X=x}[Z = 1] \geq \alpha. \quad (142)$$

Then observe that

$$\begin{aligned} & \sum_{x \in A} P_X(x) Q_{Y|X=x}[Z = 1] \\ & \geq \sum_{x \in A} P_X(x) \beta_{P_{Y|X=x}[Z=1]}(P_{Y|X=x}, Q_{Y|X=x}) \quad (143) \end{aligned}$$

$$= \sum_{x \in A} P_X(x) \beta_{P_{Y|X=x}[Z=1]} \quad (144)$$

$$\geq \beta_{P[Z=1]} \quad (145)$$

$$\geq \beta_\alpha \quad (146)$$

where (144) follows from the assumption, (146) follows because  $\beta_\alpha$  is a nondecreasing function of  $\alpha$ , and (145) is by Jensen's inequality which is applicable since  $\beta_\alpha$  is convex. Therefore, from (146) we obtain that

$$\beta_\alpha(P_X P_{Y|X}, P_X Q_{Y|X}) \geq \beta_\alpha(P_{Y|X=x}, Q_{Y|X=x}) \quad (147)$$

■ and together with (141) this concludes the proof. ■

### F. General Converse: Maximal Probability of Error

The minimax problem in (134) is generally hard to solve. A weaker bound is given by Theorem 31 which is a corollary to the next analog of Theorem 26.

*Theorem 30:* For a given code (possibly with a randomized decoder) with codewords belonging to  $\mathbb{F}$ , let

$$\begin{aligned} \epsilon &= \text{maximal error probability with } P_{Y|X} \\ \epsilon' &= \text{maximal error probability with } Q_{Y|X}. \end{aligned}$$

Then

$$\inf_{x \in \mathbb{F}} \beta_{1-\epsilon}(P_{Y|X=x}, Q_{Y|X=x}) \leq 1 - \epsilon'. \quad (148)$$

*Proof:* Consider an  $(M, \epsilon)$ -code with codewords  $\{c_j \in \mathbb{F}\}_{j=1}^M$  and a randomized decoding rule  $P_{Z|Y} : \mathbb{B} \mapsto \{0, \dots, M\}$ . We have for some  $j^*$

$$\sum_{b \in \mathbb{B}} P_{Z|Y}(j^*|b) Q_{Y|X}(b|j^*) = 1 - \epsilon' \quad (149)$$

and at the same time

$$\sum_{b \in \mathbb{B}} P_{Z|Y}(j^*|b) P_{Y|X}(b|j^*) \geq 1 - \epsilon. \quad (150)$$

Consider the hypothesis test between  $P_{Y|X=j^*}$  and  $Q_{Y|X=j^*}$  that decides in favor of  $P_{Y|X=j^*}$  only when the decoder output is  $j^*$ . By (150) the probability of correct decision under  $P_{Y|X=j^*}$  is at least  $1 - \epsilon$ , and, therefore

$$1 - \epsilon' \geq \beta_{1-\epsilon}(P_{Y|X=j^*}, Q_{Y|X=j^*}) \quad (151)$$

$$\geq \inf_{x \in \mathbb{F}} \beta_{1-\epsilon}(P_{Y|X=x}, Q_{Y|X=x}). \quad (152)$$

■

*Theorem 31 (Converse):* Every  $(M, \epsilon)$  code (maximal probability of error) with codewords belonging to  $\mathbb{F}$  satisfies

$$M \leq \inf_{Q_Y} \sup_{x \in \mathbb{F}} \frac{1}{\beta_{1-\epsilon}(x, Q_Y)} \quad (153)$$

where the infimum is over all distributions  $Q_Y$  on  $\mathbb{B}$ .

*Proof:* Repeat the argument of the proof of Theorem 27 replacing Theorem 26 by Theorem 30. ■

### G. Relation to Classical Converse Bounds

We illustrate how Theorems 26 and 30 can be used to prove all the converse results cited in Section II:

- Fano's inequality (Theorem 8): Particularize (135) to the case  $Q_Y = P_Y$ , where  $P_Y$  is the output distribution induced by the code and the channel  $P_{Y|X}$ . Note that any hypothesis test is a (randomized) binary-output transformation and therefore, by the data-processing inequality for divergence we have

$$d(1-\epsilon \mid \beta_{1-\epsilon}(P_{XY}, P_X \times P_Y)) \leq D(P_{XY} \parallel P_X \times P_Y) \quad (154)$$

where the binary divergence function satisfies

$$d(a \parallel b) = a \log \frac{a}{b} + (1-a) \log \frac{1-a}{1-b} \quad (155)$$

$$\geq -h(a) + a \log \frac{1}{b}. \quad (156)$$

Using (155) in (154), we obtain

$$\log \frac{1}{\beta_{1-\epsilon}(P_{XY}, P_X \times P_Y)} \leq \frac{I(X; Y) + h(\epsilon)}{1-\epsilon}. \quad (157)$$

Fano's inequality (19) follows from (157) and (135).

- Information spectrum converse (Theorem 11): Replace (157) with (102), which together with (135) yields

$$\frac{1}{M} \geq \beta_{1-\epsilon}(P_{XY}, P_X \times P_Y) \quad (158)$$

$$\geq \sup_{\gamma > 0} \frac{1}{\gamma} (\mathbb{P}[i(X; Y) < \log \gamma] - \epsilon). \quad (159)$$

The bound (159) is equivalent to the converse bound (23). Similarly, by using a stronger bound in place of (102) we can derive [26]. Furthermore, by keeping the freedom in choosing  $Q_Y$  in (135) we can prove a stronger version of the result.

- Wolfowitz's strong converse (Theorem 9): To apply Theorem 31 we must compute a lower bound on  $\inf_{x \in A} \beta_\alpha(x, Q_Y)$ ; but this simply amounts to taking the infimum over  $x \in A$  in (106). Thus

$$\begin{aligned} & \inf_{x \in A} \beta_\alpha(x, P_Y) \\ & \geq \sup_{\gamma > 0} \frac{1}{\gamma} \left( \alpha - \sup_{x \in A} P_{Y|X=x} \left[ \frac{dP_{Y|X=x}}{dQ_Y} \geq \gamma \right] \right). \end{aligned} \quad (160)$$

Now, suppose that  $Q_Y = P_Y$ ; then using (4) we conclude that Theorem 31 implies Theorem 9.

- Shannon–Gallager–Berlekamp (Theorem 12): Applying Theorem 31, we may first split the input space  $A$  into regions  $F_i$  such that  $\beta_\alpha(x, Q_Y)$  is constant within  $F_i$ . For example, for symmetric channels and  $Q_Y$  equal to the capacity achieving output distribution, there is no need to split  $A$  since  $\beta_\alpha(x, Q_Y)$  is identical for all  $x \in A$ . For a general DMC, we apply Theorem 26 with  $Q_{Y|X}$  chosen as follows. The distribution  $Q_{Y|X=x^n}$  only depends on the type of  $x^n$  and is chosen optimally for each type (and depending on the coding rate). Over the  $Q$ -channel, the decoder can at most distinguish codewords belonging to different types and therefore, we can estimate  $1 - \epsilon' \leq \frac{n^{|A|-1}}{M}$ . Using this estimate in (128), the proof of Theorem 12 follows along the same lines as the proof of [36, Th. 19] by weakening (128) using Chernoff-type estimates.
- Refinements to [8, Th. 12] and [11]: As we explained above, Theorem 12 is obtained from Theorem 31 by choosing  $Q_Y$  judiciously and by performing a large deviation analysis of  $\beta_\alpha$ . [8] improved Theorem 12 by extending the results to the case of infinite  $|\mathbb{B}|$  and by tightening the Chernoff-type estimates of [27]. A further improvement was found in [11] for the special case of input-symmetric channels by directly lower-bounding the average probability of error and avoiding the step of splitting a code into constant composition subcodes. Theorem 28 is tighter than the bound in [11] because for symmetric channels and relevant distributions  $Q_Y$  the value of  $\beta_\alpha(x, Q_Y)$  does not depend on  $x$  and, therefore, average probability of error is bounded directly.

### H. BSC

This section illustrates the application of the finite-length upper and lower bounds to the BSC with crossover probability  $\delta < 1/2$ .

Particularizing Theorem 15 to equiprobable input distributions and the BSC we obtain (see also [37]) the following result.

*Theorem 32:* For the BSC with crossover probability  $\delta$ , we have (161), shown at the bottom of the page.

Note that the exact evaluation of (161) poses considerable difficulties unless the blocklength is small. The next result gives a slightly weaker, but much easier to compute, bound.

*Theorem 33:* For the BSC with crossover probability  $\delta$ , there exists an  $(n, M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq \sum_{t=0}^n \binom{n}{t} \delta^t (1-\delta)^{n-t} \min \left\{ 1, (M-1) \sum_{k=0}^t \binom{n}{k} 2^{-n} \right\}. \quad (162)$$

If  $M$  is a power of 2, then the same bound holds for maximal probability of error.

*Proof:* We apply Theorem 16 (RCU bound), with  $A = B = \{0, 1\}^n$ , and the equiprobable input distribution. The information density is

$$i(x^n; y^n) = n \log(2 - 2\delta) + t \log \frac{\delta}{1 - \delta} \quad (163)$$

where  $t$  is the Hamming weight of the difference between  $x^n$  and  $y^n$ . Accordingly, since  $\bar{X}^n$  is equiprobable and independent of  $X^n$  we obtain

$$\begin{aligned} \mathbb{P} [i(\bar{X}^n; Y^n) \geq i(X^n; Y^n) \mid X^n = x^n, Y^n = y^n] \\ = \sum_{k=0}^t \binom{n}{k} 2^{-n}. \end{aligned} \quad (164)$$

The statement about the maximal probability of error is explained in Appendix A.  $\blacksquare$

It turns out that Poltyrev's bound (11), derived using linear codes and weight spectra, is in fact equal to (162) with  $M - 1$  replaced by  $2^k$ . Indeed, notice that

$$\sum_{w=0}^n \binom{n}{w} \sum_{w/2 \leq t \leq \min\{\ell, w\}} \binom{w}{t} \binom{n-w}{\ell-t} = \binom{n}{\ell} \sum_{s=0}^{\ell} \binom{n}{s}. \quad (165)$$

This holds since on the left we have counted all the ways of choosing two binary  $n$ -vectors  $X$  and  $Z$  such that  $\text{wt}(Z) = \ell$  and  $Z$  overlaps at least a half of  $X$ . The last condition is equivalent to requiring  $\text{wt}(X - Z) \leq \text{wt}(Z)$ . So we can choose  $Z$  in

$\binom{n}{\ell}$  ways and  $X$  in  $\sum_{s=0}^{\ell} \binom{n}{s}$  ways, which is the RHS of (165). Now applying (165) to (11) yields (162) with  $M - 1$  replaced by  $2^k$ .

*Theorem 34:* For the BSC with crossover probability  $\delta$ , there exists an  $(n, M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq \sum_{t=0}^n \binom{n}{t} \min \{ \delta^t (1 - \delta)^{n-t}, (M - 1) 2^{-n-1} \}. \quad (166)$$

If  $M$  is a power of 2, then the same bound holds for maximal probability of error.

*Proof:* Taking  $P_{X^n}$  to be equiprobable on  $\{0, 1\}^n$ , the DT bound of Theorem 17 is equal to  $\frac{M+1}{2}$  times the minimal probability of error of an optimal binary hypothesis test between  $n$  fair coin tosses (with prior probability  $\frac{M-1}{M+1}$ ) and  $n$  bias- $\delta$  coin tosses (with prior probability  $\frac{2}{M+1}$ ). The upper bound (67) on the average error probability becomes

$$\epsilon \leq \mathbb{E} \left[ 2^{-[na - bZ - \log \frac{M-1}{2}]^+} \right] \quad (167)$$

where

$$a = 1 + \log_2(1 - \delta) \quad (168)$$

$$b = \log_2 \frac{1 - \delta}{\delta} \quad (169)$$

and  $Z \sim B(n, \delta)$  is a binomial random variable with parameters  $n$  and  $\delta$ . Averaging over  $Z$ , (167) becomes (166). The statement about the maximal probability of error is explained in Appendix A.  $\blacksquare$

For comparison, Feinstein's lemma (Theorem 1), with equiprobable input distribution yields

$$M^*(n, \epsilon) \geq \sup_{t > 0} 2^{nt} (\epsilon - \mathbb{P}[Z \geq n(a - t)/b]) \quad (170)$$

where  $Z \sim B(n, \delta)$ .

Gallager's random coding bound (7) also with equiprobable input distribution yields<sup>9</sup>

$$\log_2 M^*(n, \epsilon) \geq n E_r^{-1} \left( \frac{1}{n} \log_2 \frac{1}{\epsilon} \right) \quad (171)$$

where [16, Th. 5.6.2, Cor. 2 and Example 1 in Sec. 5.6.]

$$E_r(1 - h(s)) = \begin{cases} d(s|\delta) & s \in (\delta, s^*] \\ h(s) - 2 \log_2 s_1 & s > s^* \end{cases} \quad (172)$$

<sup>9</sup>Inequality (171) holds for average probability of error. Figs. 1 and 2 show the corresponding bound on maximal error probability where we drop the half of the codewords with worse error probability. This results in an additional term of  $-1$  appended to the RHS of (171), while  $\frac{1}{\epsilon}$  becomes  $\frac{2}{\epsilon}$  therein.

$$\mathbb{E} [\epsilon(X_1, \dots, X_M)] =$$

$$= 1 - 2^{n-nM} \sum_{i=0}^n \binom{n}{i} \delta^i (1 - \delta)^{n-i} \sum_{\ell=0}^{M-1} \binom{n}{i}^{\ell} \frac{1}{1 + \ell} \binom{M-1}{\ell} \left( \sum_{j=i+1}^n \binom{n}{j} \right)^{M-1-\ell}. \quad (161)$$

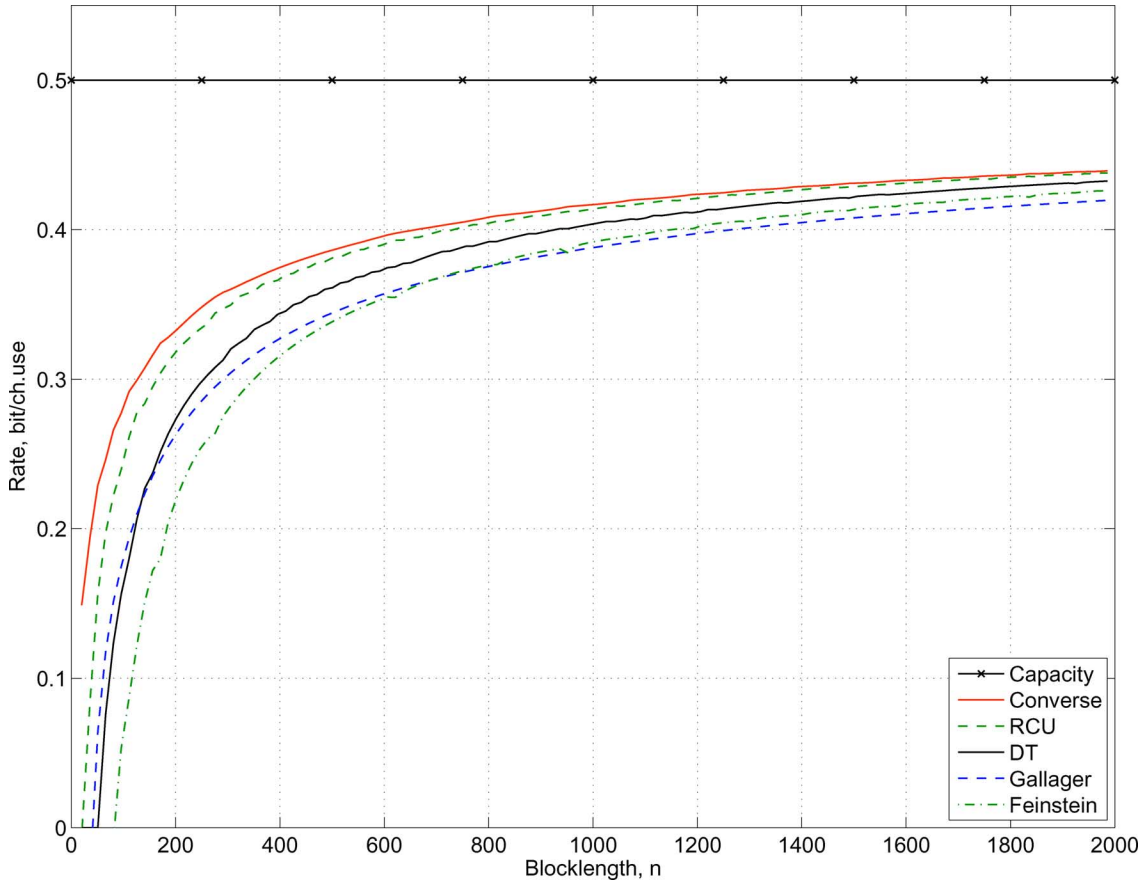


Fig. 1. Rate-blocklength tradeoff for the BSC with crossover probability  $\delta = 0.11$  and maximal block error rate  $\epsilon = 10^{-3}$ .

and  $s^* = \frac{\sqrt{\delta}}{\sqrt{\delta} + \sqrt{1-\delta}}$ ,  $s_1 = \sqrt{\delta} + \sqrt{1-\delta}$ .

We now turn our attention to the computation of the converse bound of Theorem 28. Choosing  $Q_{Y^n}$  equiprobable on  $\{0, 1\}^n$  we recover the classical sphere packing bound (cf. [16, eq. (5.8.19)] for an alternative expression).

*Theorem 35:* For the BSC with crossover probability  $\delta$ , the size of an  $(n, M, \epsilon)$  code (average error probability) must satisfy

$$M \leq \frac{1}{\beta_{1-\epsilon}^n}. \quad (173)$$

where  $\beta_\alpha^n$  is defined as

$$\beta_\alpha^n = (1 - \lambda)\beta_L + \lambda\beta_{L+1} \quad (174)$$

with

$$\beta_\ell = \sum_{k=0}^{\ell} \binom{n}{k} 2^{-n} \quad (175)$$

where  $0 \leq \lambda < 1$  and the integer  $L$  are defined by

$$\alpha = (1 - \lambda)\alpha_L + \lambda\alpha_{L+1} \quad (176)$$

with

$$\alpha_\ell = \sum_{k=0}^{\ell-1} \binom{n}{k} (1 - \delta)^{n-k} \delta^k. \quad (177)$$

*Proof:* To streamline notation, we denote  $\beta_\alpha^n = \beta_\alpha(x^n, Q_{Y^n})$  since it does not depend on  $x^n$ , and  $Q_{Y^n}$  is fixed. Then, the Hamming weight of the output word is a sufficient statistic for discriminating between  $P_{Y^n|X^n=0}$  and  $Q_{Y^n}$ . Thus, the optimal randomized test is

$$P_{Z_0|Y^n}(1|y^n) = \begin{cases} 0, & |y^n| > L_\alpha^n, \\ \lambda_\alpha^n, & |y^n| = L_\alpha^n, \\ 1, & |y^n| < L_\alpha^n, \end{cases} \quad (178)$$

where  $L_\alpha^n \in \mathbb{Z}_+$  and  $\lambda_\alpha^n \in [0, 1)$  are uniquely determined by

$$\sum_{y^n \in \mathcal{A}} P_{Y^n|X^n}(y^n|\mathbf{0}) P_{Z_0|Y^n}(1|y^n) = \alpha. \quad (179)$$

Then we find that

$$\beta_\alpha^n = \lambda_\alpha^n \binom{n}{L_\alpha^n} 2^{-n} + \sum_{k=0}^{L_\alpha^n-1} \binom{n}{k} 2^{-n}. \quad (180)$$

Thus, by Theorem 28

$$M^*(n, \epsilon) \leq \frac{1}{\beta_{1-\epsilon}^n}. \quad (181)$$

■

The numerical evaluation of (162), (166), and (173) is shown in Figs. 1 and 2, along with the bounds by Feinstein (170) and Gallager (171). As we anticipated analytically, the DT bound is always tighter than Feinstein's bound. For  $\delta = 0.11$  and

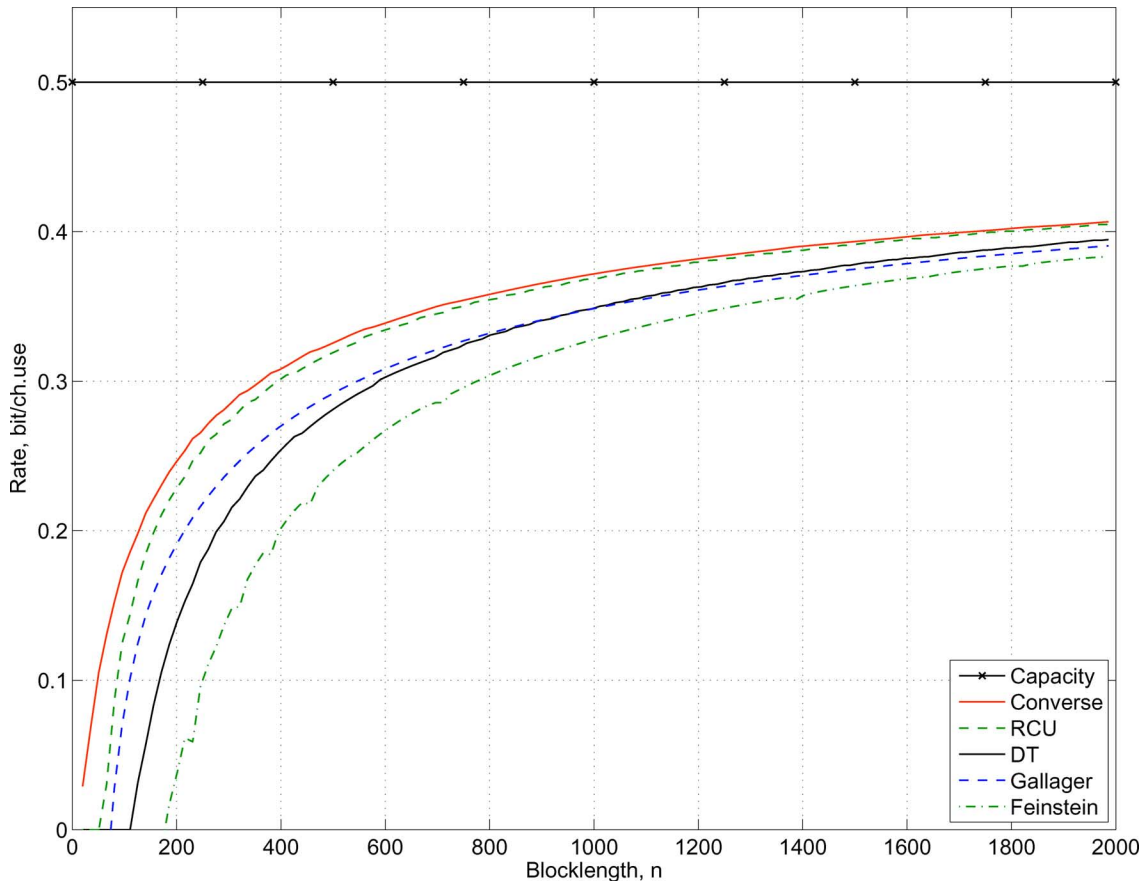


Fig. 2. Rate-blocklength tradeoff for the BSC with crossover probability  $\delta = 0.11$  and maximal block error rate  $\epsilon = 10^{-6}$ .

$\epsilon = 0.001$ , we can see in Fig. 1 that for blocklengths greater than about 150, Theorem 17 gives better results than Gallager's bound. In fact, for large  $n$  the gap to the converse upper bound of the new lower bound is less than half that of Gallager's bound. The RCU achievability bound (162) is uniformly better than all other bounds. In fact for all  $n \geq 20$  the difference between (162) and the converse is within 3–4 bits in  $\log_2 M$ . This tendency remains for other choices of  $\delta$  and  $\epsilon$ , although, for smaller  $\epsilon$  and/or  $\delta$ , Gallager's bound (originally devised to analyze the regime of exponentially small  $\epsilon$ ) is tighter for a larger range of blocklengths, see Fig. 2. A similar relationship between the three bounds holds, qualitatively, in the case of the additive white Gaussian noise channel (Section III-J).

### I. BEC

Next we illustrate the application of the achievability bounds in Theorems 15, 16, 17, and 22 to the special case of the binary erasure channel. Using Theorem 15 we obtain the next bound.

*Theorem 36:* For the BEC with erasure probability  $\delta$ , we have (182), shown at the bottom of the page.

Easier to evaluate is the DT bound (Theorem 17), which particularizes to the following.

*Theorem 37:* For the BEC with erasure probability  $\delta$ , there exists an  $(n, M, \epsilon)$  code (average probability of error) such that

$$\epsilon \leq \sum_{t=0}^n \binom{n}{t} \delta^t (1-\delta)^{n-t} 2^{-[n-t-\log_2(\frac{M-1}{2})]^+}. \quad (183)$$

If  $M$  is a power of 2, then the same bound holds for maximal probability of error. In any case there exists an  $(n, M, \epsilon)$  code (maximal probability of error) such that

$$\epsilon \leq \sum_{t=0}^n \binom{n}{t} \delta^t (1-\delta)^{n-t} 2^{-[n-t-\log_2(M-1)]^+}. \quad (184)$$

*Proof:* Using Theorem 17 with  $A = \{0, 1\}^n$ ,  $B = \{0, e, 1\}^n$  and the equiprobable input distribution, it

$$\mathbb{E}[\epsilon(X_1, \dots, X_M)] = 1 - \sum_{\ell=0}^{M-1} \binom{M-1}{\ell} \frac{1}{\ell+1} \sum_{j=0}^n \binom{n}{j} \delta^{n-j} (1-\delta)^j 2^{-j\ell} (1-2^{-j})^{M-1-\ell}. \quad (182)$$



follows that if  $y^n$  contains  $t$  erasures and coincides with  $x^n$  in all the nonerased bits, then

$$i(x^n; y^n) = n - t \quad (185)$$

and otherwise,  $i(x^n; y^n) = -\infty$ . Then (67) implies (183) since  $t$  erasures happen with probability  $\binom{n}{t} \delta^t (1 - \delta)^{n-t}$ . If  $M$  is a power of 2 then the same bound holds for maximal probability of error by using linear codes (see Appendix A). Bound (184) is obtained by exactly the same argument, except that Theorem 22 must be used in lieu of Theorem 17. ■

Application of Theorem 16 yields exactly (184) but only for average probability of error. Since Theorem 16 is always stronger than Gallager's bound, we conclude that Theorem 37 is also stronger than Gallager's bound for the BEC and therefore achieves the random coding exponent. Similarly Theorem 21 (and hence Theorem 22) is always stronger than Feinstein's bound, see (92). Therefore, Theorem 37 is also stronger than Feinstein's bound for the BEC. The average block erasure probability for a random ensemble of linear codes is given in [38]; it can be shown that it is sandwiched between (183) and (184), which are also considerably easier to compute.

Upper bounds on error probability can be converted easily into lower bounds on  $M^*(n, \epsilon)$ . For example, Theorem 37 for the maximal probability of error formalism implies<sup>10</sup>

$$M^*(n, \epsilon) \geq \max \left\{ 2^k : \mathbb{E} \left[ 2^{-\lfloor Z - \log_2 \frac{2^k - 1}{2} \rfloor^+} \right] \leq \epsilon \right\} \quad (186)$$

where  $Z \sim B(n, \delta)$  is a binomial random variable with parameters  $n$  and  $\delta$ .

The upper bound on code size given by Theorem 31 (with capacity achieving output distribution) is improved by the following result,<sup>11</sup> which is stronger than related bounds such as in [39].

**Theorem 38:** For the BEC with erasure probability  $\delta$ , the average error probability of an  $(n, M, \epsilon)$  code satisfies

$$\epsilon \geq \sum_{\ell=\lfloor n - \log_2 M \rfloor + 1}^n \binom{n}{\ell} \delta^\ell (1 - \delta)^{n-\ell} \left( 1 - \frac{2^{n-\ell}}{M} \right) \quad (187)$$

even if the encoder knows the location of the erasures non-causally.

*Proof:* It is easy to show that the probability of correct decoding in an  $M$ -ary equiprobable hypothesis testing problem where the observable takes one out of  $J$  values is upper bounded by  $J/M$ , even if stochastic decision rules are allowed. Indeed,

<sup>10</sup>For numerical purposes we can safely weaken (186) by replacing  $\log_2 \frac{2^k - 1}{2}$  with  $(k - 1)$ .

<sup>11</sup>For a  $q$ -ary erasure channel, Theorem 38 holds replacing  $2^{n-\ell-k}$  by  $q^{n-\ell-k}$  and  $\log_2$  by  $\log_q$ . In fact, this  $q$ -ary analog of (187) is achievable by  $q$ -ary maximum distance separable (MDS) codes.

suppose that the true hypothesis is a (random variable)  $A$ , the observable output is  $B$  and the decision is  $C$ ; then

$$\mathbb{P}[C = A] = \sum_{a=1}^M \sum_{b=1}^J P_A(a) P_{B|A}(b|a) P_{C|B}(a|b) \quad (188)$$

$$= \sum_{b=1}^J \frac{1}{M} \sum_{a=1}^M P_{C|B}(a|b) P_{B|A}(b|a) \quad (189)$$

$$\leq \sum_{b=1}^J \frac{1}{M} \sum_{a=1}^M P_{C|B}(a|b) \quad (190)$$

$$\leq \frac{J}{M}. \quad (191)$$

Now suppose that the location of the erasures is known to the encoder, and there are  $z \in \{0, \dots, n\}$  erasures. Then, regardless of the code (possibly dependent on the erasure pattern) chosen by the encoder, the decoder faces an  $M$ -ary equiprobable hypothesis testing problem where the observable takes one out of  $2^{n-z}$  values. Therefore, the probability of error is lower bounded by  $\left[ 1 - \frac{2^{n-z}}{M} \right]^+$ . Since each pattern of  $z$  erasures occurs with probability  $(1 - \delta)^{n-z} \delta^z$  and there are  $\binom{n}{z}$  of them, (187) follows. ■

Figs. 3 and 4 show that, as expected, (183) is quite a bit tighter than the Gallager and Feinstein bounds. In fact, the gap between (183) and (187) is below 3 bits in  $\log_2 M$ , uniformly across the blocklengths shown on the plot. Fig. 5 compares the DT bound (183) with the BEC achievability bound (12); they are within one bit of each other, the winner depending on a particular value of  $n$ . The zigzagging of the plot of (12) is a behavior common to all bounds that are restricted to integer values of  $\log_2 M$ . The complexity of the computation of (12) is  $O(n^3)$ , compared to  $O(n)$  for the DT bound (183).

## J. The AWGN Channel

**1) The Channel and Power Constraints:** For the real-valued additive-noise white Gaussian channel we have the following specific definitions:

- $\mathbf{A} = \mathbb{R}^n$ ,
- $\mathbf{B} = \mathbb{R}^n$  and
- $P_{Y^n|X^n=x^n} = \mathcal{N}(x^n, \mathbf{I}_n)$ .

Additionally, codewords are subject to one of three types of power constraints:

- equal-power constraint:  $M_e^*(n, \epsilon, P)$  denotes the maximal number of codewords, such that each codeword  $c_i \in X^n$  satisfies

$$\|c_i\|^2 = nP. \quad (192)$$

- maximal power constraint:  $M_m^*(n, \epsilon, P)$  denotes the maximal number of codewords, such that each codeword  $c_i \in X^n$  satisfies

$$\|c_i\|^2 \leq nP. \quad (193)$$

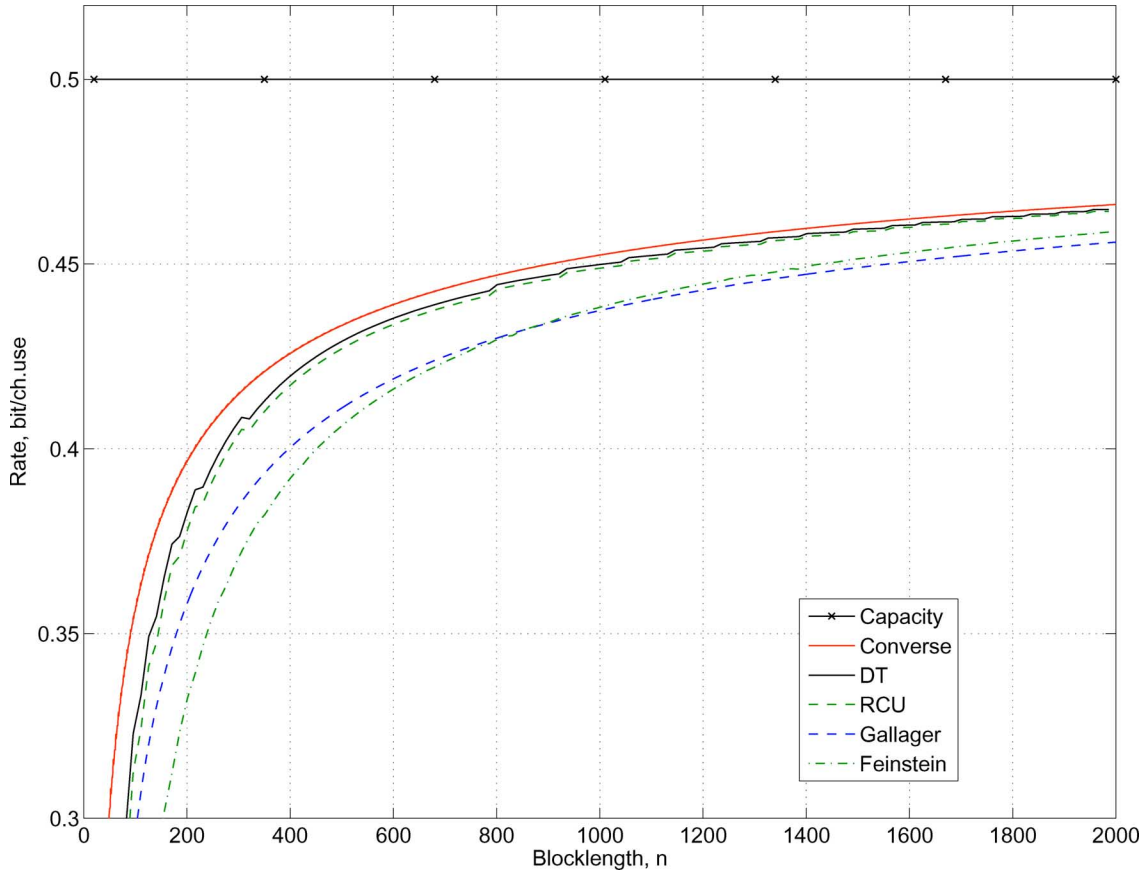


Fig. 3. Rate-blocklength tradeoff for the BEC with erasure probability  $\delta = 0.5$  and maximal block error rate  $\epsilon = 10^{-3}$ .

- average power constraint:  $M_a^*(n, \epsilon, P)$  denotes the maximal size  $M$  of a codebook that satisfies

$$\frac{1}{M} \sum_{i=1}^M \|c_i\|^2 \leq nP. \tag{194}$$

It is easiest to analyze  $M_e^*$ , but  $M_m^*$  and  $M_a^*$  are more interesting from the practical viewpoint. Following Shannon [3] we will make use of simple inequalities relating all three quantities, summarized in the following.

*Lemma 39:* For any  $0 < P < P'$  the inequalities

$$M_e^*(n, \epsilon, P) \leq M_m^*(n, \epsilon, P) \leq M_e^*(n + 1, \epsilon, P) \tag{195}$$

and

$$M_m^*(n, \epsilon, P) \leq M_a^*(n, \epsilon, P) \leq \frac{1}{1 - P/P'} M_m^*(n, \epsilon, P') \tag{196}$$

hold.

*Proof:* The left-hand bounds are obvious. The right-hand bound in (195) follows from the fact that we can always take the  $M_m^*$ -code and add an  $(n + 1)$ -th coordinate to each codeword to equalize the total power to  $nP$ . The right-hand bound in (196) is a consequence of the Chebyshev inequality on the probability of finding a codeword with power greater than  $nP'$  in the  $M_a^*$ -code. ■

The particularization of the exact error probability achieved by random coding in Theorem 15 leads to (41) which turns out to be the tightest of all the bounds for the AWGN channel. However the particularization of the  $\kappa\beta$ -bound to the AWGN channel is of paramount importance in Section IV.

2) *Evaluation of  $\beta_\alpha^n$ :* We will now apply Theorems 25 and 28 to the AWGN channel with equal-power constraint (192). For each  $n$ , the set  $F_n$  of permissible inputs is

$$F_n \triangleq \{x^n : \|x^n\|^2 = nP\} \subset \mathbb{R}^n. \tag{197}$$

To use Theorems 25 and 28 we must also choose the auxiliary distribution  $P_{Y^n}$  over  $\mathcal{B}^n$ . A particularly convenient choice is

$$P_{Y^n} = \mathcal{N}(0, \sigma_Y^2 \mathbf{I}_n). \tag{198}$$

with  $\sigma_Y^2$  to be specified later. Due to the spherical symmetry of both  $F_n$  and (198), for all  $x \in F_n$

$$\beta_\alpha^n(x, \mathcal{N}(0, \sigma_Y^2 \mathbf{I}_n)) = \beta_\alpha^n. \tag{199}$$

To simplify calculations, we choose  $x = x_0 = (\sqrt{P}, \sqrt{P}, \dots, \sqrt{P})$ . The information density is given by

$$i(x_0; y^n) = \frac{n}{2} \log \sigma_Y^2 + \frac{\log e}{2} \sum_{i=1}^n \left[ \frac{y_i^2}{\sigma_Y^2} - (y_i - \sqrt{P})^2 \right]. \tag{200}$$

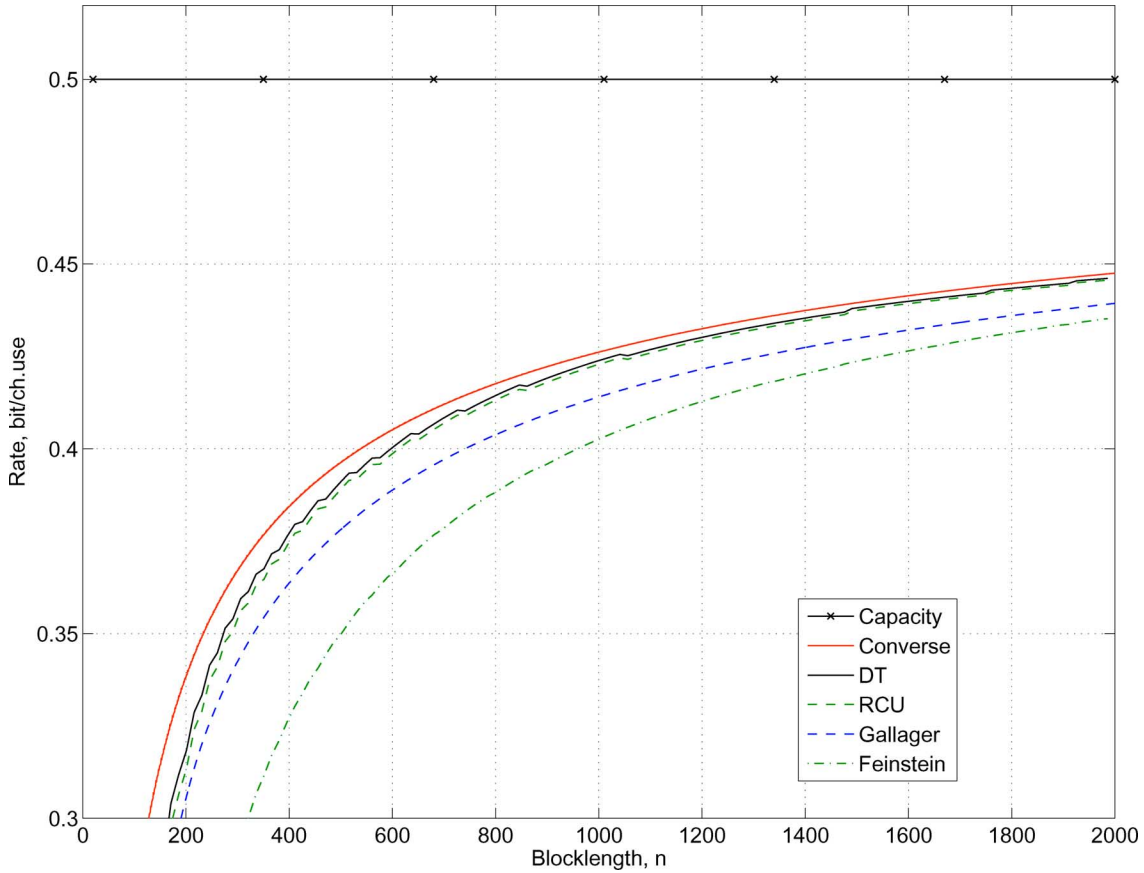


Fig. 4. Rate-blocklength tradeoff for the BEC with erasure probability  $\delta = 0.5$  and maximal block error rate  $\epsilon = 10^{-6}$ .

It is convenient to define independent standard Gaussian variables  $Z_i \sim \mathcal{N}(0, 1)$ ,  $i = 1, \dots, n$ . Then, under  $P_{Y^n}$  and under  $P_{Y^n|X^n=x_0}$ , the information density  $i(x_0; Y^n)$  has the same distribution as

$$G_n = n \log \sigma_Y - n \frac{P}{2} \log e + \frac{1}{2} \log e \sum_{i=1}^n \left( (1 - \sigma_Y^2) Z_i^2 + 2\sqrt{P} \sigma_Y Z_i \right) \quad (201)$$

and

$$H_n = n \log \sigma_Y + n \frac{P}{2\sigma_Y^2} \log e + \frac{1}{2\sigma_Y^2} \log e \sum_{i=1}^n \left( (1 - \sigma_Y^2) Z_i^2 + 2\sqrt{P} Z_i \right) \quad (202)$$

respectively. A judicious choice is

$$\sigma_Y^2 = 1 + P \quad (203)$$

since it maximizes  $D(P_{Y^n|X^n=x_0} || P_{Y^n}) = \mathbb{E}[H_n]$ , and  $P_{Y^n}$  coincides with the capacity-achieving output distribution for the AWGN channel. With this choice of  $\sigma_Y^2$ , (201) and (202) become

$$G_n = \frac{n}{2} \log(1 + P) - \frac{P}{2} \sum_{i=1}^n \left( 1 + Z_i^2 - 2\sqrt{1 + \frac{1}{P}} Z_i \right) \log e \quad (204)$$

and

$$H_n = \frac{n}{2} \log(1 + P) + \frac{1}{2} \frac{P}{(1 + P)} \sum_{i=1}^n \left( 1 - Z_i^2 + \frac{2}{\sqrt{P}} Z_i \right) \log e. \quad (205)$$

Finally, using the Neyman–Pearson lemma (Appendix B), we obtain the following result.

*Theorem 40:* For the additive white Gaussian noise channel and all  $x \in \mathbb{F}_n$

$$\beta_\alpha^n = \beta_\alpha^n(x, \mathcal{N}(0, \sigma_Y^2 \mathbf{I}_n)) = \mathbb{P}[G_n \geq \gamma] \quad (206)$$

where  $\gamma$  satisfies

$$\mathbb{P}[H_n \geq \gamma] = \alpha. \quad (207)$$

Applying Theorems 28, 40, and Lemma 39, we obtain the following converse bound.

*Theorem 41:* For the AWGN channel and for any  $n$  and  $\epsilon$  (average probability of error) we have

$$M_m^*(n - 1, \epsilon, P) \leq M_e^*(n, \epsilon, P) \leq \frac{1}{\mathbb{P}[G_n \geq \gamma_n]} \quad (208)$$

where  $\gamma_n$  satisfies

$$\mathbb{P}[H_n \geq \gamma_n] = 1 - \epsilon \quad (209)$$

and  $G_n$  and  $H_n$  are defined in (204) and (205).

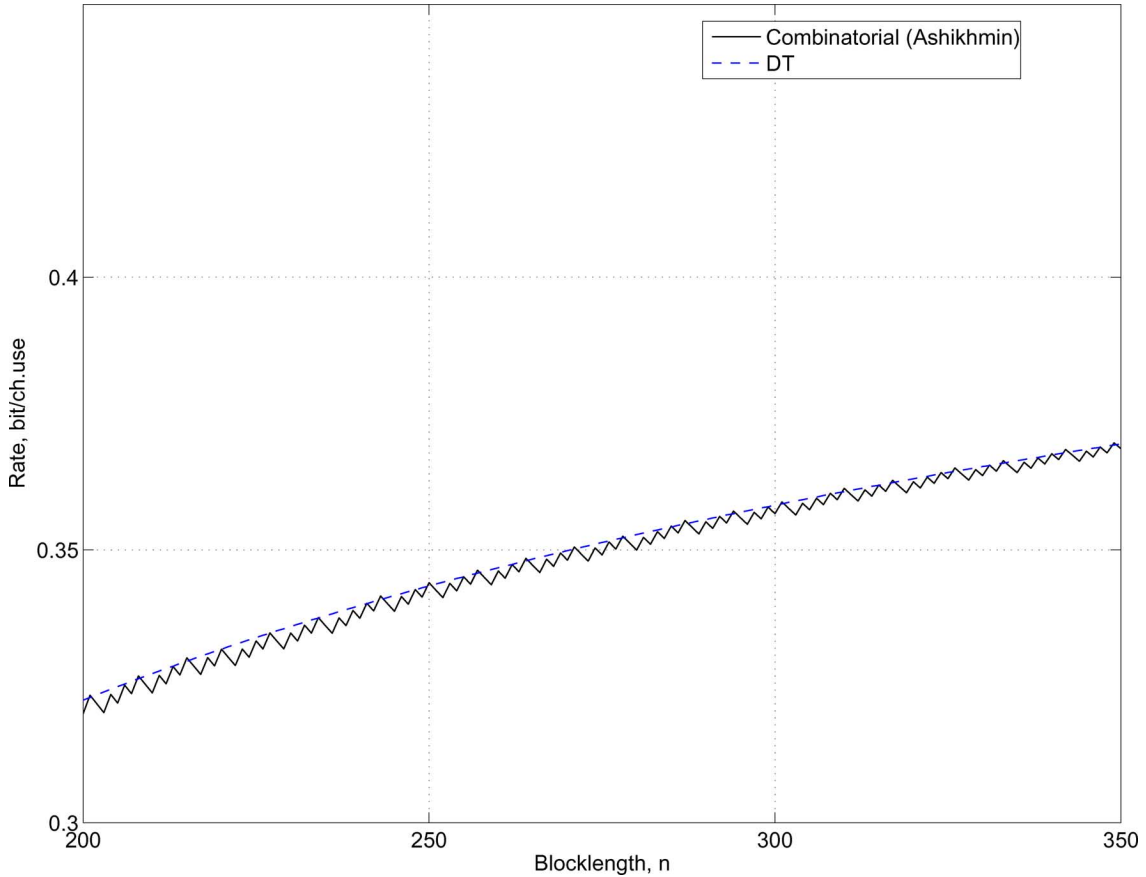


Fig. 5. Comparison of the DT-bound (183) and the combinatorial bound of Ashikhmin (12) for the BEC with erasure probability  $\delta = 0.5$  and probability of block error  $\epsilon = 10^{-3}$ .

The distributions of  $G_n$  and  $H_n$  are noncentral  $\chi^2$ . However, the value of  $\mathbb{P}[G_n \geq \gamma]$  decreases exponentially, and for large  $n$ , traditional series expansions of the noncentral  $\chi^2$  distribution do not work very well; a number of other techniques must be used to evaluate these probabilities, including Chernoff bounding as well as (106) and (103).

3) *Evaluation of  $\kappa_\tau^n$* : Although we are free to choose any  $P_{Y^n}$ , it is convenient to use (198).

*Theorem 42*: For the chosen  $P_{Y^n}$ ,  $F_n$  and for any  $\tau \in [0, 1]$  and  $n \geq 1$ , we have

$$\kappa_\tau^n(F_n, P_{Y^n}) = P_0 \left[ \frac{p_1(r)}{p_0(r)} \geq \gamma \right] \quad (210)$$

where  $\gamma$  satisfies

$$P_1 \left[ \frac{p_1(r)}{p_0(r)} \geq \gamma \right] = \tau \quad (211)$$

with  $p_0$  and  $p_1$  being probability density functions (PDFs) of  $P_0$  and  $P_1$ , defined as

$$p_0(r) = \frac{r^{n/2-1} e^{-r/(2+2P)}}{(2+2P)^{n/2} \Gamma(n/2)} \quad (212)$$

$$p_1(r) = \frac{1}{2} e^{-(r+nP)/2} \left( \frac{r}{nP} \right)^{n/4-1/2} I_{n/2-1}(\sqrt{nP}r) \quad (213)$$

where  $I_a(y)$  is a modified Bessel function of a first kind:

$$I_a(y) = (y/2)^a \sum_{j=0}^{\infty} \frac{(y^2/4)^j}{j! \Gamma(a+j+1)}. \quad (214)$$

The proof is given in Appendix D.

A straightforward application of a (local) central limit theorem yields the following result.

*Lemma 43*: Under the conditions of Theorem 42

$$\lim_{n \rightarrow \infty} \kappa_\tau^n = \kappa_\tau^\infty = 1 - 2Q(r_\tau^*) \quad (215)$$

where

$$r_\tau^* = \frac{\sqrt{1+2P}}{1+P} Q^{-1} \left( \frac{1-\tau}{2} \right). \quad (216)$$

Experimentally, we have observed that the convergence in (215) is very fast. For example, for  $P = 1$ ,  $n = 10$  and  $\tau \in [10^{-6}, 10^{-1}]$ , we find that

$$|\kappa_\tau^n - \kappa_\tau^\infty| \leq 5 \cdot 10^{-3} \kappa_n(\tau). \quad (217)$$

Summarizing, we have particularized Theorems 25 and 31 to the AWGN channel to show

$$\frac{\kappa_\tau^n}{\beta_{1-\epsilon+\tau}^n} \leq M_m^*(n, \epsilon, P) \leq \frac{1}{\beta_{1-\epsilon}^{n+1}} \quad (218)$$

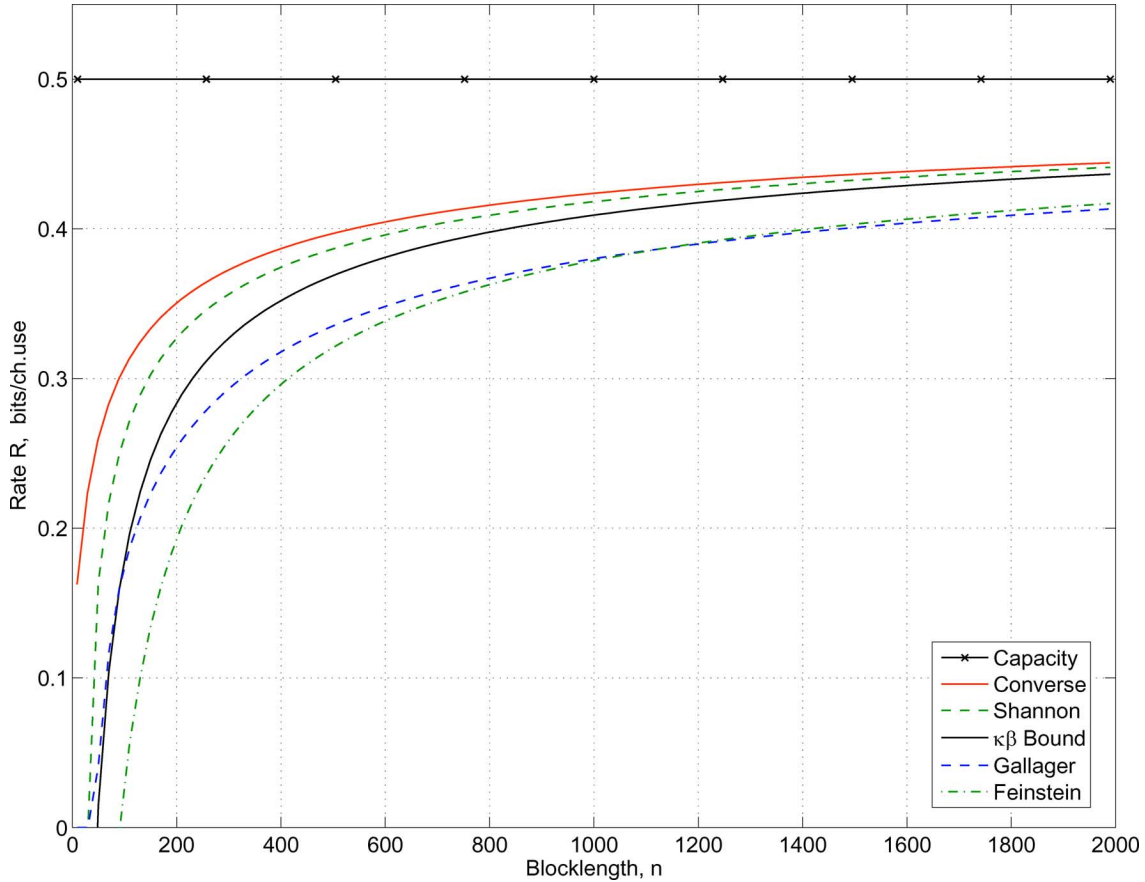


Fig. 6. Bounds for the AWGN channel,  $SNR = 0$  dB,  $\epsilon = 10^{-3}$ .

where  $\beta_\alpha^n$  and  $\kappa_\tau^n$  determined by Theorems 40 and 42.

4) *Numerical Evaluation:* In this section our goal is to compare various achievability bounds. To emphasize the quality of the bounds we also compare them against the converse, Theorem 41. As usual, we plot the converse bounds for the average probability of error formalism and achievability bounds for the maximal probability of error formalism. The power constraint is the maximal one, i.e., we are plotting bounds on  $M_m^*(n, \epsilon)$ . The results are found on Figs. 6 and 7. Let us first explain how each bound was computed:

- 1) The converse bound is Theorem 41. Note that in [3] Shannon gives another converse bound (38). However, in this case both bounds numerically coincide almost exactly and for this reason only the bound in Theorem 41 is plotted.
- 2) Feinstein’s bound is Theorem 24 with

$$F = \{x^n : \|x^n\|^2 \leq nP\} \tag{219}$$

and  $P_{X^n} = \mathcal{N}(0, P\mathbf{I}_n)$ .

- 3) Gallager’s bound is Theorem 14, where we optimize the choice of  $\delta$  for each  $R$ , and then select the largest  $R$  that still keeps the bound (44) below the required  $\epsilon$ .
- 4) The  $\kappa_\beta$  bound is an application of Theorem 25 with  $\beta_\alpha$  and  $\kappa_\tau$  given by Theorems 40 and 42. As discussed earlier, the convergence in (215) is very fast and  $\kappa_\tau^n$  affects rate only as  $\frac{\log \kappa_\tau^n}{n}$ ; thus we can safely replace the  $\kappa_\tau^n$  with  $\kappa_\tau^\infty$ . In this way, for each  $n$  we need to compute only  $\beta_{1-\epsilon}^n$ .

- 5) Shannon’s bound<sup>12</sup>: The bound in (41) is on *average* probability of error. For the BSC and BEC we transformed from average to maximal probability of error using the random linear code method. Unfortunately, for the AWGN channel we could not find anything equivalent; instead we need to recourse to traditional “purging”. Namely, if we have an  $(M, \tau\epsilon)$ -code for average probability then there must exist a  $(\tau M, \epsilon)$ -subcode for maximal probability. Consequently, if  $M_S(n, \epsilon)$  is the maximal cardinality of the codebook guaranteed by the Shannon bound, then instead we plot

$$M_S^{max}(n, \epsilon) = \max_{\tau \in [0,1]} (1 - \tau)M_S(n, \tau\epsilon). \tag{220}$$

Shannon’s achievability bound is the clear winner on both Figs. 6 and 7. It comes very close to the converse; for example, on Fig. 6 in terms of  $\log_2 M$  the difference between the Shannon bound and the converse is less than 6 bits uniformly across the range of blocklengths depicted on the graph. This illustrates that random codes are not only optimal asymptotically, but also almost optimal even for rather small blocklengths.

The drawback of the Shannon bound is that it is harder to compute and analyze than the  $\kappa_\beta$  bound and requires a “purging” procedure to guarantee a small maximal probability of error. Section IV-C invokes the  $\kappa_\beta$  bound to analyze the

<sup>12</sup>We use expression (42) and the representation of  $q_n(\phi)$  as a noncentral  $t$ -distribution given by [3, (17)]. Note that to improve numerical stability of the integration in (42) it is convenient to multiply the integrand by  $(\sin \theta(M))^{-(n-2)}$ .

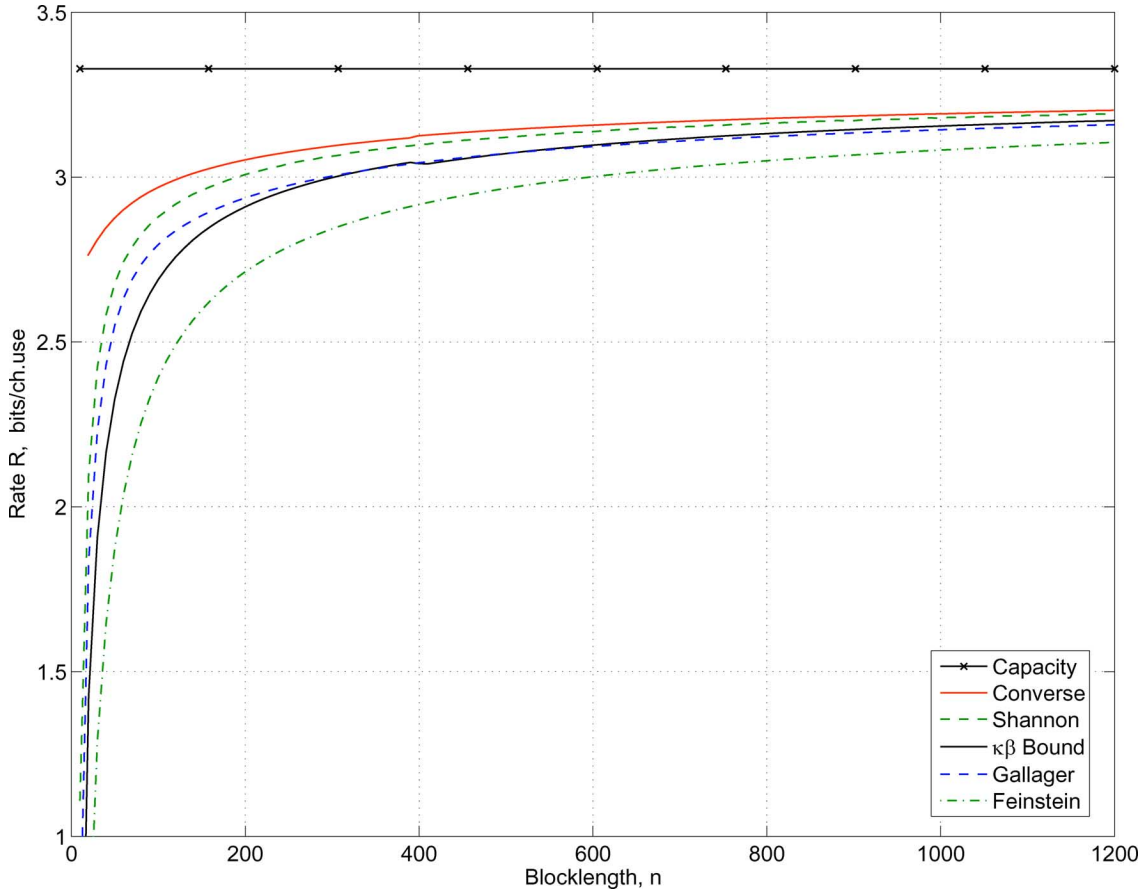


Fig. 7. Bounds for the AWGN channel,  $SNR = 20$  dB,  $\epsilon = 10^{-6}$ .

asymptotic expansion of  $\log M_m^*(n, \epsilon)$ . In Figs. 6 and 7 we can see that the  $\kappa\beta$  bound is also quite competitive for finite  $n$ .

Comparing the  $\kappa\beta$  bound and the classical bounds of Feinstein and Gallager, we see that, as expected, the  $\kappa\beta$  bound is uniformly better than Feinstein's bound. In the setup of Fig. 6, the  $\kappa\beta$  bound is a significant improvement over Gallager's bound, coming very close to the Shannon bound as well as the converse. In Fig. 7, both the  $\kappa\beta$  and Gallager bounds are again very close to the Shannon bound but this time Gallager's bound is better for small  $n$ . There are two reasons for this. First, recall that we have analyzed a suboptimal decoder based on hypothesis testing, whereas Gallager used the maximum likelihood decoder. It seems that for small  $n$  it is important to use optimal decoding. Moreover, Gallager's analysis is targeted at very small  $\epsilon$ . Indeed, as we go from  $10^{-3}$  to  $10^{-6}$ , the tightness of Gallager's bound improves significantly. In general we observe that Gallager's bound improves as the channel becomes better and as  $\epsilon$  gets smaller. On the other hand, the  $\kappa\beta$  bound is much more uniform over both SNR and  $\epsilon$ . In Section IV, the  $\kappa\beta$  bound, in contrast to Gallager's bound, yields the correct  $\sqrt{n}$  term in the asymptotic expansion of  $\log M_m^*(n, \epsilon)$ .

Comparing the RCU bound and the DT bound (and its relative, the  $\kappa\beta$  bound), the DT bound is very handy theoretically and does not lose much nonasymptotically compared to the RCU bound. In fact, for the BEC the DT bound is tighter than the RCU bound. Also, the DT bound (in the form of Theorems 22 and 25) and the  $\kappa\beta$  bound are directly applicable to the

maximal probability of error, whereas the RCU bound requires further manipulation (e.g., Appendix A).

#### IV. NORMAL APPROXIMATION

We turn to the asymptotic analysis of the maximum achievable rate for a given blocklength. In this section, our goal is to show a normal-approximation refinement of the coding theorem. To that end, we introduce the following definition.

*Definition 1:* The channel dispersion  $V$  (measured in squared information units per channel use) of a channel with capacity  $C$  is equal to

$$V = \lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \left( \frac{nC - \log M^*(n, \epsilon)}{Q^{-1}(\epsilon)} \right)^2 \quad (221)$$

$$= \lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \frac{(nC - \log M^*(n, \epsilon))^2}{2 \ln \frac{1}{\epsilon}}. \quad (222)$$

In fact, we show that for both discrete memoryless channels and Gaussian channels,

$$\log M^*(n, \epsilon) = nC - \sqrt{nV} Q^{-1}(\epsilon) + O(\log n). \quad (223)$$

The asymptotic behavior in (223) is particularly useful in conjunction with the nonasymptotic upper and lower bounds developed in Section III, as (223) turns out to be an accurate and succinct approximation to the fundamental finite blocklength limit for even rather short blocklengths and rates well below capacity.

Thus, an excellent approximation to the rate penalty incurred for operating at blocklength  $n$  and error probability  $\epsilon$  is

$$\frac{\log M^*(n, \epsilon)}{n} \approx C - A(n, \epsilon) \quad (224)$$

where  $A(n, \epsilon)$  is the  $A$  required for the probability of error of the binary equiprobable hypothesis test

$$H_0 : Z_i = A + N_i \quad i = 1, \dots, n \quad (225)$$

$$H_1 : Z_i = -A + N_i \quad i = 1, \dots, n \quad (226)$$

to be  $\epsilon$  if  $\{N_i\}$  are independent and Gaussian with variances  $V$ . This implies that if the target is to transmit at a given fraction of capacity  $0 < \eta < 1$  and at a given error probability  $\epsilon$ , the required blocklength scales linearly with the channel dispersion:

$$n^*(\eta, \epsilon) \approx \frac{V}{C^2} \left( \frac{Q^{-1}(\epsilon)}{1 - \eta} \right)^2. \quad (227)$$

An important tool in this section is the following nonasymptotic result.

*Theorem 44 (Berry–Esseen):* (e.g., Theorem 2, [40, Ch. XVI.5]) Let  $X_k$ ,  $k = 1, \dots, n$  be independent with

$$\mu_k = \mathbb{E}[X_k] \quad (228)$$

$$\sigma_k^2 = \text{Var}[X_k] \quad (229)$$

$$t_k = \mathbb{E}[|X_k - \mu_k|^3] \quad (230)$$

$$\sigma^2 = \sum_{k=1}^n \sigma_k^2 \quad (231)$$

$$T = \sum_{k=1}^n t_k. \quad (232)$$

Then for any<sup>13</sup> $-\infty < \lambda < \infty$

$$\left| \mathbb{P} \left[ \sum_{k=1}^n (X_k - \mu_k) \geq \lambda \sigma \right] - Q(\lambda) \right| \leq \frac{6T}{\sigma^3}. \quad (233)$$

#### A. DMC

The DMC has finite input alphabet  $\mathcal{A}$ , finite output alphabet  $\mathcal{B}$ , and conditional probabilities

$$P_{Y^n|X^n}(y^n|x^n) = \prod_{i=1}^n W(y_i|x_i) \quad (234)$$

where  $W(\cdot|x)$  is a conditional probability mass function on  $\mathcal{B}$  for all  $x \in \mathcal{A}$ , which is abbreviated as  $W_x$  when notationally convenient. We denote the simplex of probability distributions on  $\mathcal{A}$ , by  $\mathcal{P}$ . It is useful to partition  $\mathcal{P}$  into  $n$ -types

$$\mathcal{P}_n = \{P \in \mathcal{P} : nP(x) \in \mathbb{Z}_+ \ \forall x \in \mathcal{A}\}. \quad (235)$$

<sup>13</sup>Note that for i.i.d.  $X_k$  it is known [41] that the factor of 6 in (233) can be replaced by 0.7975. In this paper, the exact value of the constant does not affect the results and so we take the conservative value of 6 even in the i.i.d. case.

We denote by  $M^*(n, \epsilon)$  (respectively,  $M_{\text{avg}}^*(n, \epsilon)$ ) the cardinality of the largest codebook with maximal (respectively, average) probability of error below  $\epsilon$ . We use the following notation and terminology:

- divergence variance

$$V(P||Q) = \sum_{x \in \mathcal{A}} P(x) \left[ \log \frac{P(x)}{Q(x)} \right]^2 - D(P||Q)^2. \quad (236)$$

- conditional divergence variance

$$V(W||Q|P) = \sum_{x \in \mathcal{A}} P(x) V(W_x||Q). \quad (237)$$

- output distribution  $PW$  as

$$PW(y) = \sum_{x \in \mathcal{A}} P(x) W(y|x).$$

- mutual information

$$I(P, W) = \mathbb{E}[i(X; Y)] = \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P(x) W(y|x) \log \frac{W(y|x)}{PW(y)}. \quad (238)$$

- unconditional information variance

$$U(P, W) = \text{Var}(i(X; Y)) \quad (239)$$

$$= \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P(x) W(y|x) \log^2 \frac{W(y|x)}{PW(y)} - [I(P, W)]^2 \quad (240)$$

$$= V(P \times W || P \times PW) \quad (241)$$

- conditional information variance

$$V(P, W) = \mathbb{E}[\text{Var}(i(X; Y) | X)] \quad (242)$$

$$= \sum_{x \in \mathcal{A}} P(x) \left\{ \sum_{y \in \mathcal{B}} W(y|x) \log^2 \frac{W(y|x)}{PW(y)} - [D(W_x || PW)]^2 \right\} \quad (243)$$

$$= V(W || PW | P) \quad (244)$$

- third absolute moment of the information density

$$T(P, W) = \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P(x) W(y|x) \times \left| \log \frac{W(y|x)}{PW(y)} - D(W_x || PW) \right|^3. \quad (245)$$

Note that  $V(W||Q|P)$  is defined only provided that  $W_x \ll Q$  for  $P$ -almost all  $x$ , and the divergence variance is defined only if  $P \ll Q$ . Continuity of  $U(P, W)$ ,  $V(P, W)$  and  $T(P, W)$  is established by Lemma 62 in Appendix E.

The compact subset of *capacity-achieving distributions*  $\Pi$  is

$$\Pi \triangleq \{P \in \mathcal{P} : I(P, W) = C\} \quad (246)$$

where

$$C = \max_{P \in \mathcal{P}} I(P, W). \quad (247)$$

1) *Achievability Bound:*

*Theorem 45:* For any  $P \in \mathcal{P}$ , we have

$$\log M_{\text{avg}}^*(n, \epsilon) \geq nI(P, W) - \sqrt{nU(P, W)} Q^{-1}(\epsilon) + O(1) \quad (248)$$

$$\log M^*(n, \epsilon) \geq nI(P, W) - \sqrt{nU(P, W)} Q^{-1}(\epsilon) - \frac{1}{2} \log n + O(1) \quad (249)$$

if  $U(P, W) > 0$  and

$$\log M^*(n, \epsilon) \geq nI(P, W) + \log \epsilon \quad (250)$$

if  $U(P, W) = 0$ . Finally, if  $V(W_x \| PW) > 0$  whenever  $P(x) > 0$  then

$$\log M^*(n, \epsilon) \geq nI(P, W) - \sqrt{nU(P, W)} Q^{-1}(\epsilon) + O(1). \quad (251)$$

*Proof:* Select  $P \in \mathcal{P}$ . Let  $\mathbf{A} = \mathcal{A}^n$ , and choose the product measure  $P^n$  as the distribution of  $X^n$ . Passing this distribution through  $W^n$  induces a joint probability distribution on  $(X^n, Y^n)$ , and the information density is the sum of independent identically distributed  $Z_k$

$$i(X^n; Y^n) = \sum_{k=1}^n \log \frac{W(Y_k | X_k)}{PW(Y_k)} = \sum_{k=1}^n Z_k. \quad (252)$$

The random variable  $Z_k$  has the distribution of  $i(X; Y)$  when  $(X, Y)$  is distributed according to  $P \times W$ . Accordingly, it has mean  $I(P, W)$  and variance  $U(P, W)$ , and its third absolute moment is bounded according to the following auxiliary result whose proof is in Appendix F.

*Lemma 46:*

$$\mathbb{E} [|i(X; Y) - I(X; Y)|^3] \leq \left( (|\mathcal{A}|^{1/3} + |\mathcal{B}|^{1/3}) \times 3e^{-1} \log e + \log \min\{|\mathcal{A}|, |\mathcal{B}|\} \right)^3. \quad (253)$$

Suppose that  $U(P, W) = 0$ , and therefore  $i(X^n; Y^n) = nI(P, W)$ . Taking  $\log \beta = nI(P, W) - \delta$  for an arbitrary  $\delta > 0$  in Theorem 1 we get (250).

Now, assume that  $U(P, W) > 0$  and denote

$$B \triangleq \frac{6\kappa}{U(P, W)^{3/2}} \quad (254)$$

where  $\kappa$  is the RHS of (253).

To use the DT bound (67) we need to prove that for some  $\gamma$  the following inequality holds:

$$\epsilon \geq \mathbb{E} \left[ \exp \left\{ -[i(X^n; Y^n) - \log \gamma]^+ \right\} \right] \quad (255)$$

$$= \mathbb{P}[i(X^n; Y^n) \leq \log \gamma] \quad (256)$$

$$+ \gamma \mathbb{E} \left[ \exp \left\{ -i(X^n; Y^n) \right\} 1_{\{i(X^n; Y^n) > \log \gamma\}} \right]. \quad (257)$$

Denote for arbitrary  $\tau$

$$\log \gamma = nI(P, W) - \tau \sqrt{nU(P, W)}. \quad (258)$$

According to Theorem 44, we have

$$|\mathbb{P}[i(X^n; Y^n) \leq \log \gamma] - Q(\tau)| \leq \frac{B}{\sqrt{n}}. \quad (259)$$

For sufficiently large  $n$ , let

$$\tau = Q^{-1} \left( \epsilon - \left( \frac{2 \log 2}{\sqrt{2\pi}} + 5B \right) \frac{1}{\sqrt{n}} \right). \quad (260)$$

Then, from (259) we obtain

$$\mathbb{P}[i(X^n; Y^n) \leq \log \gamma] \leq \epsilon - 2 \left( \frac{\log 2}{\sqrt{2\pi}} + 2B \right) \frac{1}{\sqrt{n}}. \quad (261)$$

We now bound the second term (257) by the following technical result proved in Appendix G.

*Lemma 47:* Let  $Z_1, Z_2, \dots, Z_n$  be independent random variables,  $\sigma^2 = \sum_{j=1}^n \text{Var} Z_j$  be nonzero and  $T = \sum_{j=1}^n \mathbb{E} [|Z_j - \mathbb{E} Z_j|^3] < \infty$ ; then for any  $A$

$$\mathbb{E} \left[ \exp \left\{ -\sum_{j=1}^n Z_j \right\} 1_{\left\{ \sum_{j=1}^n Z_j > A \right\}} \right] \leq 2 \left( \frac{\log 2}{\sqrt{2\pi}} + \frac{12T}{\sigma^2} \right) \frac{1}{\sigma} \exp\{-A\}. \quad (262)$$

Therefore, we have

$$\gamma \mathbb{E} \left[ \exp \left\{ -i(X^n; Y^n) \right\} 1_{\{i(X^n; Y^n) > \log \gamma\}} \right] \leq 2 \left( \frac{\log 2}{\sqrt{2\pi}} + 2B \right) \frac{1}{\sqrt{n}}. \quad (263)$$

Summing (261) and (263) we prove inequality (255). Hence, by Theorem 17 we get

$$\log M^*(n, \epsilon) \geq \log \gamma \quad (264)$$

$$= nI - \tau \sqrt{nU} \quad (265)$$

$$= nI - \sqrt{nU(P, W)} Q^{-1}(\epsilon) + O(1) \quad (266)$$

because according to (260) and the differentiability of  $Q^{-1}$  we have

$$\tau = Q^{-1}(\epsilon) + O \left( \frac{1}{\sqrt{n}} \right). \quad (267)$$

Note that (248) implies (249) after applying

$$M^*(n, \epsilon) \geq (1 - \tau_n) M_{\text{avg}}^*(n, \epsilon \tau_n) \quad (268)$$

with  $\tau_n = 1 - \frac{1}{\sqrt{n}}$ .

Finally, the proof of (251) repeats the proof of (248) step-by-step with the only change that Theorem 21 is used



instead of Theorem 17 and in (254) we replace  $U(P, W)$  by  $\min_x V(W_x || PW)$ . ■

Note that by using the Feinstein bound (5) we could only prove (249), not the stronger (248) or (251). This suboptimality in the  $\log n$  term is an analytical expression of the fact that we have already observed in Section III: namely, that the Feinstein bound is not tight enough for the refined analysis of  $\log M^*(n, \epsilon)$ .

As another remark, we recall that by using the DT bound, Theorem 17, we proved that with input distribution  $P$  we can select  $\exp\{nI(P, W) - \sqrt{nU(P, W)}Q^{-1}(\epsilon)\}$  messages which are distinguishable with probability of error  $\epsilon$ . It is not hard to see that by using the  $\kappa\beta$  bound, Theorem 25, we could select<sup>14</sup>  $\exp\{nI(P, W) - \sqrt{nV(P, W)}Q^{-1}(\epsilon)\}$ , which is the same for a capacity achieving  $P$  (see Lemma 62) and is larger otherwise. While in the unconstrained case we used the DT bound, in the cost constrained cases we resort to the  $\kappa\beta$  bound (as in the AWGN case treated in Section IV-C).

*Converse Theorem for DMC:* We need to define a few new quantities in order to state the converse counterpart to Theorem 45.

- Define maximal and minimal conditional variances (they exist since  $V(P, W)$  is continuous) as

$$V_{\max} = \max_{P \in \Pi} V(P, W) = \max_{P \in \Pi} U(P, W) \quad (269)$$

and

$$V_{\min} = \min_{P \in \Pi} V(P, W) = \min_{P \in \Pi} U(P, W). \quad (270)$$

- Define the (unique) capacity achieving output distribution  $P_Y^*$  by  $P_Y^* = P^*W$ , where  $P^*$  is any capacity achieving input distribution.
- $W$  is an *exotic DMC* if  $V_{\max} = 0$  and there exists an input letter  $x_0$  such that: a) for any capacity achieving  $P$ :  $P(x_0) = 0$ , b)  $D(W_{x_0} || P_Y^*) = C$ , and c)  $V(W_{x_0} || P_Y^*) > 0$ . (See Appendix H for an example of an exotic DMC.)
- For any  $P_0 \in \mathcal{P}_n$  denote a *type* of elements  $x^n \in \mathcal{A}^n$  by

$$T_{P_0}^n \triangleq \{x^n : \forall a \in \mathcal{A} : \sum_{i=1}^n 1_{\{x_i=a\}} = P_0(a)\}. \quad (271)$$

- For any  $n$  and  $P_0 \in \mathcal{P}_n$  denote by  $M_{P_0}^*(n, \epsilon)$  the maximal cardinality of the codebook with codewords in  $T_{P_0}^n$  and maximal probability of error below  $\epsilon$ .

*Theorem 48:* Fix a DMC  $W$ .

- If  $0 < \epsilon \leq \frac{1}{2}$ , then there exists a constant  $F > 0$  such that for all  $P_0 \in \mathcal{P}_n$  and all sufficiently large  $n$

$$\log M_{P_0}^*(n, \epsilon) \leq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + \frac{1}{2} \log n + F. \quad (272)$$

- If  $\frac{1}{2} < \epsilon < 1$  and the DMC is not exotic, then there exists a constant  $F > 0$  such that for all  $P_0 \in \mathcal{P}_n$  and all sufficiently large  $n$

$$\log M_{P_0}^*(n, \epsilon) \leq nC - \sqrt{nV_{\max}}Q^{-1}(\epsilon) + \frac{1}{2} \log n + F. \quad (273)$$

<sup>14</sup>Theorem 25 is applied with  $Q_Y = (PW)^n$  and  $\mathbf{F} = T_{P^*}^n$ , a  $P$ -type in the input space  $\mathcal{A}^n$ . The analysis of  $\beta$  is the same as in the proof of Theorem 48, Section IV-B.II; for  $\kappa$  it is sufficient to use the lower bound (121).

- If  $\frac{1}{2} < \epsilon < 1$  and the DMC is exotic, then there exists a constant  $G > 0$  such that for all  $P_0 \in \mathcal{P}_n$  and all sufficiently large  $n$

$$\log M_{P_0}^*(n, \epsilon) \leq nC + Gn^{1/3}. \quad (274)$$

*Proof:* See Appendix I ■

2) *DMC Dispersion:* The following result is a refinement of [31].

*Theorem 49:* For a DMC and  $0 < \epsilon \leq 1/2$  we have

$$\log M^*(n, \epsilon) = nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + O(\log n) \quad (275)$$

where  $C$  is the capacity and  $V_{\min}$  is the minimal variance of the information density over all capacity achieving distributions (cf. (270)). In addition, if there exists a capacity achieving input distribution  $P^*$  such that  $V(W_x || P^*W) > 0$  whenever  $P^*(x) > 0$  then

$$\log M^*(n, \epsilon) \geq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + O(1). \quad (276)$$

*Proof:* Theorem 45 yields, by taking  $P \in \mathcal{P}$  to be a distribution that achieves capacity and minimizes  $U(P, W)$  (or  $V(P, W)$  since they coincide on  $\Pi$  by Lemma 62),

$$\log M^*(n, \epsilon) \geq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + O(\log n). \quad (277)$$

For the lower bound, take  $n \geq N_0$  for  $N_0$  from Theorem 48. Then any  $(n, M, \epsilon)$  is composed of subcodes over types  $T_{P_0}^n$  for  $P_0 \in \mathcal{P}_n$ . If we remove all codewords except those in  $T_{P_0}^n$  and leave the decoding regions untouched, then we obtain an  $(n, M'_{P_0}, \epsilon)$  code over  $T_{P_0}^n$ . But then Theorem 48 states that

$$\begin{aligned} \log M'_{P_0} &\leq \log M_{P_0}^*(n, \epsilon) \\ &\leq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + \frac{1}{2} \log n + F. \end{aligned} \quad (278)$$

Since  $M$  is a sum of  $M'_{P_0}$  over all  $P_0 \in \mathcal{P}_n$  and the cardinality of  $\mathcal{P}_n$  is no more than  $(n+1)^{|\mathcal{A}|-1}$ , we conclude

$$\begin{aligned} \log M^*(n, \epsilon) &\leq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + \left(|\mathcal{A}| - \frac{1}{2}\right) \log n + F'. \end{aligned} \quad (279)$$

This completes the proof of (275) and (276) follows from (251). ■

It is useful to introduce the following definition.

*Definition 2:* For a channel with  $\epsilon$ -capacity  $C_\epsilon$ , the  $\epsilon$ -dispersion is defined for  $\epsilon \in (0, 1) - \{\frac{1}{2}\}$  as

$$V_\epsilon = \limsup_{n \rightarrow \infty} \frac{1}{n} \left( \frac{nC_\epsilon - \log M^*(n, \epsilon)}{Q^{-1}(\epsilon)} \right)^2. \quad (280)$$

Note that for  $\epsilon < \frac{1}{2}$ , approximating  $\frac{1}{n} \log M^*(n, \epsilon)$  by  $C_\epsilon$  is optimistic and smaller dispersion is preferable, while for  $\epsilon > \frac{1}{2}$ , it is pessimistic and larger dispersion is more favorable. Since  $Q^{-1}(\frac{1}{2}) = 0$ , it is immaterial how to define  $V_{\frac{1}{2}}$  as far as the normal approximation (223) is concerned.

Invoking the strong converse, we can show that the  $\epsilon$ -dispersion of a DMC is

$$V_\epsilon = \begin{cases} V_{\min} & \epsilon < 1/2 \\ V_{\max} & \epsilon > 1/2. \end{cases} \quad (281)$$

Because of the importance of channel dispersion, we note the following upper bound (see also [16, Exercise 5.23]).

*Theorem 50:* For the DMC with  $\min\{|\mathcal{A}|, |\mathcal{B}|\} > 2$  we have

$$V \leq 2 \log^2 \min\{|\mathcal{A}|, |\mathcal{B}|\} - C^2. \quad (282)$$

For the DMC with  $\min\{|\mathcal{A}|, |\mathcal{B}|\} = 2$  we have

$$V \leq 1.2 \log^2 e - C^2. \quad (283)$$

*Proof:* This is a simple consequence of Lemma 62 in Appendix E. ■

Since the typical blocklength needed to achieve capacity is governed by  $V/C^2$ , it is natural to ask whether for very small capacities the upper-bound in (282) can be improved to prevent the divergence of  $\frac{V}{C^2}$ . Such a bound is not possible over all  $W$  with fixed alphabet sizes, since such a collection of DMCs always includes all of the BSCs for which we know that  $\frac{V}{C^2} \rightarrow \infty$  as  $C \rightarrow 0$ .

We briefly consider the normal approximation in the case of average error probability. Recall that  $M_{\text{avg}}^*(n, \epsilon)$  stands for the maximal cardinality of a codebook with average probability of error below  $\epsilon$ . Then, dropping all codewords whose probabilities of error are above  $\tau\epsilon$ ,  $\tau > 1$  (see the comment at the beginning of Section III-D), we obtain

$$M^*(n, \epsilon) \leq M_{\text{avg}}^*(n, \epsilon) \leq \frac{\tau}{\tau - 1} M^*(n, \tau\epsilon). \quad (284)$$

Carefully following the proof of the converse we can conclude that the  $O(\log n)$  term in the upper bound on  $\log M^*$  does not have any singularities in a neighborhood of any  $\epsilon \in (0, 1)$ . So we can claim that, for  $\tau$  sufficiently close to 1, the expansion

$$\log M^*(n, \tau\epsilon) = nC - \sqrt{nV_{\min}}Q^{-1}(\tau\epsilon) + O(\log n) \quad (285)$$

holds uniformly in  $\tau$ . Now, setting  $\tau_n = 1 + \frac{1}{\sqrt{n}}$ , we obtain

$$\log M_{\text{avg}}^*(n, \epsilon) \leq nC - \sqrt{nV_{\min}}Q^{-1}\left(\epsilon + \frac{\epsilon}{\sqrt{n}}\right) + O(\log n). \quad (286)$$

Expanding  $Q^{-1}$  by Taylor's formula and using the lower bound on  $M_{\text{avg}}^*$  in (284) we obtain the following result.

*Corollary 51:* For a DMC, if  $0 < \epsilon \leq 1/2$ , we have

$$\log M_{\text{avg}}^*(n, \epsilon) = nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + O(\log n) \quad (287)$$

$$\log M_{\text{avg}}^*(n, \epsilon) \geq nC - \sqrt{nV_{\min}}Q^{-1}(\epsilon) + O(1). \quad (288)$$

We note the following differences with Strassen's treatment of the normal approximation for DMCs in [31]. First, the DT bound allows us to prove that the  $\log n$  term cannot be negative<sup>15</sup>. Second, we streamline the proof in the case  $\epsilon < 1/2$  by

<sup>15</sup>This estimate of the  $\log n$  term cannot be improved without additional assumptions, because the BEC has zero  $\log n$  term; see Theorem 53.

using Lemma 64 to obtain the expansion. In contrast, an expansion up to the order  $o(\sqrt{n})$  can be obtained with considerably less effort by using Lemma 63. Third, [31] argues that the case  $\epsilon > 1/2$  can be treated similarly, whereas we demonstrate that this is only true for nonexotic channels as a result of the difference between using Lemma 63 and Lemma 64. (See the counter-example after the proof of Lemma 63 in Appendix J and also the discussion of exotic channels in Appendix H.) Fourth, we prove the expansion for  $\log M_{\text{avg}}^*$  (i.e., for the average probability of error formalism).

3) *Application to the BSC and the BEC:* For the BSC and BEC we can improve upon the  $O(\log n)$  term given by Theorem 49.

*Theorem 52:* For the BSC with crossover probability  $\delta$ , such that  $\delta \notin \{0, \frac{1}{2}, 1\}$ , we have

$$\log_2 M^*(n, \epsilon) = n(1 - h(\delta)) - \sqrt{n\delta(1 - \delta)} \log_2 \frac{1 - \delta}{\delta} Q^{-1}(\epsilon) + \frac{1}{2} \log_2 n + O(1), \quad (289)$$

regardless of whether  $\epsilon$  is maximal or average probability of error.

*Proof:* Appendix K. ■

Interestingly, Gallager's bound does not yield a correct  $\sqrt{n}$  term in (54); the Feinstein, DT, and RCU bounds all yield the correct  $\sqrt{n}$  term for the BSC; Feinstein's bound has worse  $\log n$  term than the DT bound. Finally, only the RCU bound (162) achieves the optimal  $\log n$  term.

*Theorem 53:* For the BEC with erasure probability  $\delta$ , we have

$$\log_2 M^*(n, \epsilon) = n(1 - \delta) - \sqrt{n\delta(1 - \delta)}Q^{-1}(\epsilon) + O(1) \quad (290)$$

regardless of whether  $\epsilon$  is maximal or average probability of error.

*Proof:* Appendix K. ■

For the BEC, Gallager's bound does not achieve the correct lower-order terms in (54); Feinstein's bound yields the correct  $\sqrt{n}$  term but a suboptimal  $\log n$  term; both DT bounds (Theorems 17 and 22) and the RCU bound achieve the optimal  $\log n$  term.

## B. The AWGN Channel

*Theorem 54:* For the AWGN channel with SNR  $P$ ,  $0 < \epsilon < 1$  and for equal-power, maximal-power and average-power constraints,

$$\log M^*(n, \epsilon, P) = nC - \sqrt{nV}Q^{-1}(\epsilon) + O(\log n) \quad (291)$$

where

$$C = \frac{1}{2} \log(1 + P) \quad (292)$$

$$V = \frac{P}{2} \frac{P + 2}{(P + 1)^2} \log^2 e. \quad (293)$$

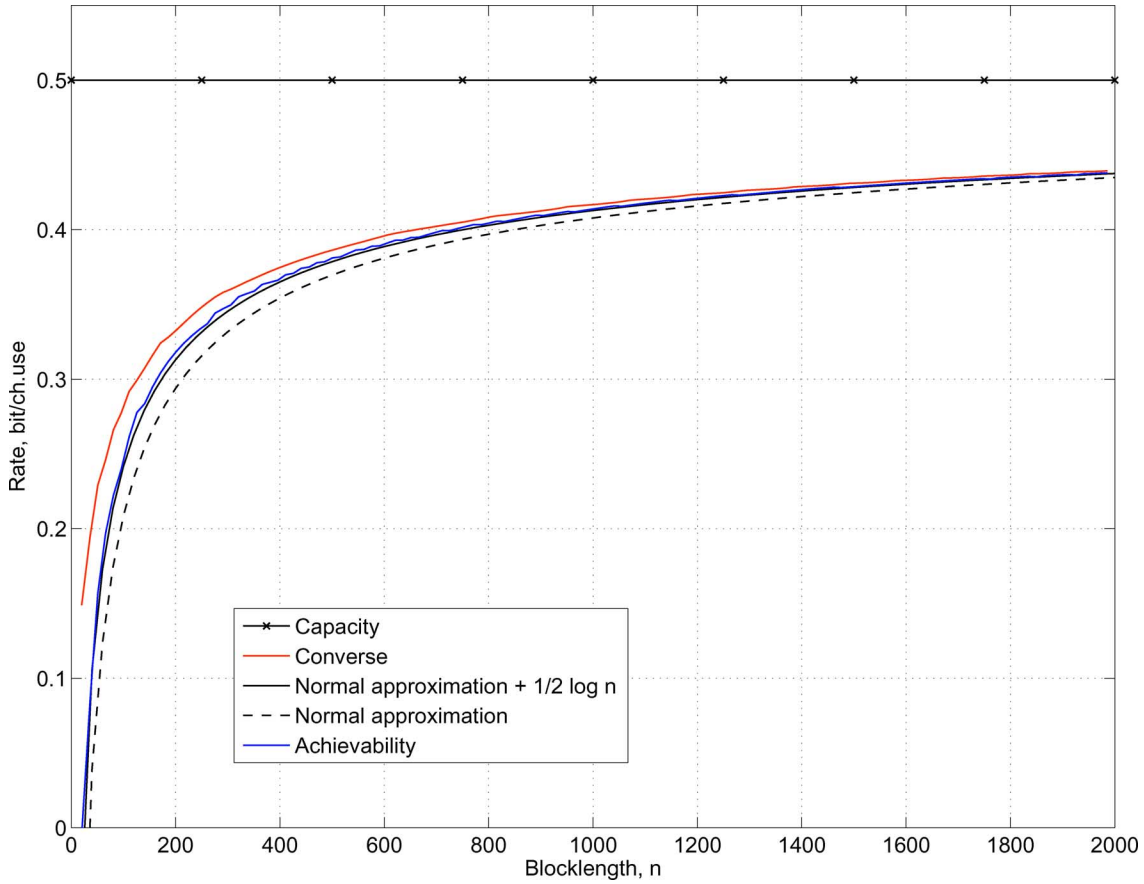


Fig. 8. Rate-blocklength tradeoff for the BSC with crossover probability  $\delta = 0.11$  and maximal block error rate  $\epsilon = 10^{-3}$ .

More precisely, for equal-power and maximal-power constraints, the  $O(\log n)$  term in (291) can be bounded by

$$\begin{aligned} O(1) &\leq \log M_{e,m}^*(n, \epsilon, P) - [nC - \sqrt{nV}Q^{-1}(\epsilon)] \\ &\leq \frac{1}{2} \log n + O(1) \end{aligned} \tag{294}$$

whereas for average-power constraint we have

$$\begin{aligned} O(1) &\leq \log M_a^*(n, \epsilon, P) - [nC - \sqrt{nV}Q^{-1}(\epsilon)] \\ &\leq \frac{3}{2} \log n + O(1). \end{aligned} \tag{295}$$

*Proof:* Appendix L. ■

The approximation in Theorem 54 (up to  $o(\sqrt{n})$ ) is attributed in [7] to Shannon [3] for the case of equipower codewords.<sup>16</sup> However, in Theorem 54 the rate is changing with  $n$ , while expressions [3, eqs. (9) and (73)] are not directly applicable here because they are asymptotic equivalence relations for fixed rate. Similarly, an asymptotic expansion up to the  $o(\sqrt{n})$  term is put forward in [47] based on a heuristic appeal to the central-limit theorem and fine quantization of the input/output alphabets.

<sup>16</sup>A different  $\frac{1}{\sqrt{n}}$  term is claimed in [7] for the case of codebook-averaged power which is not compatible with Theorem 54.

### C. Normal Approximation versus Finite Blocklength Bounds

In Figs. 8–11, we compare the normal approximation (289) and (290) to the tight bounds, computed in Section III-H (BSC) and Section III-I (BEC), correspondingly. Similarly, Figs. 12 and 13 depict the normal approximation (291) for  $\log M_m^*(n, \epsilon)$  (maximal power constraint) along with the bounds (208) and (220) for the AWGN channel. In view of (294) and the empirical evidence, we have chosen the following as a normal approximation for the AWGN channel:

$$\log M^*(n, \epsilon) \approx nC - \sqrt{nV}Q^{-1}(\epsilon) + \frac{1}{2} \log n. \tag{296}$$

Although generally pessimistic, the normal approximation is excellent for blocklengths higher than 200 (BSC(0.11) and BEC(0.5) with  $\epsilon = 10^{-3}$  and AWGN, SNR = 20 dB with  $\epsilon = 10^{-6}$ ) and 800 (AWGN, SNR = 0 dB,  $\epsilon = 10^{-3}$  and BSC(0.11),  $\epsilon = 10^{-6}$ ). The conclusion from these figures is that the normal approximation is quite accurate when transmitting at a large fraction (say  $> 0.8$ ) of channel capacity. For example, in the Table I we show the numerical results for the blocklength required by the converse, guaranteed by the achievability and predicted by error-exponents and normal approximation<sup>17</sup> for achieving rate  $R = 0.9C$ .

<sup>17</sup>For the BSC and the AWGN channel we use the approximation formula (289) which has an additional  $\frac{1}{2} \log n$  term. For the AWGN channel the DT bound is replaced by the  $\kappa\beta$  bound. The error-exponent approximation is  $N \approx -\frac{1}{E(R)} \log \frac{1}{\epsilon}$ , where  $E(R)$  is known since the rate is above critical.

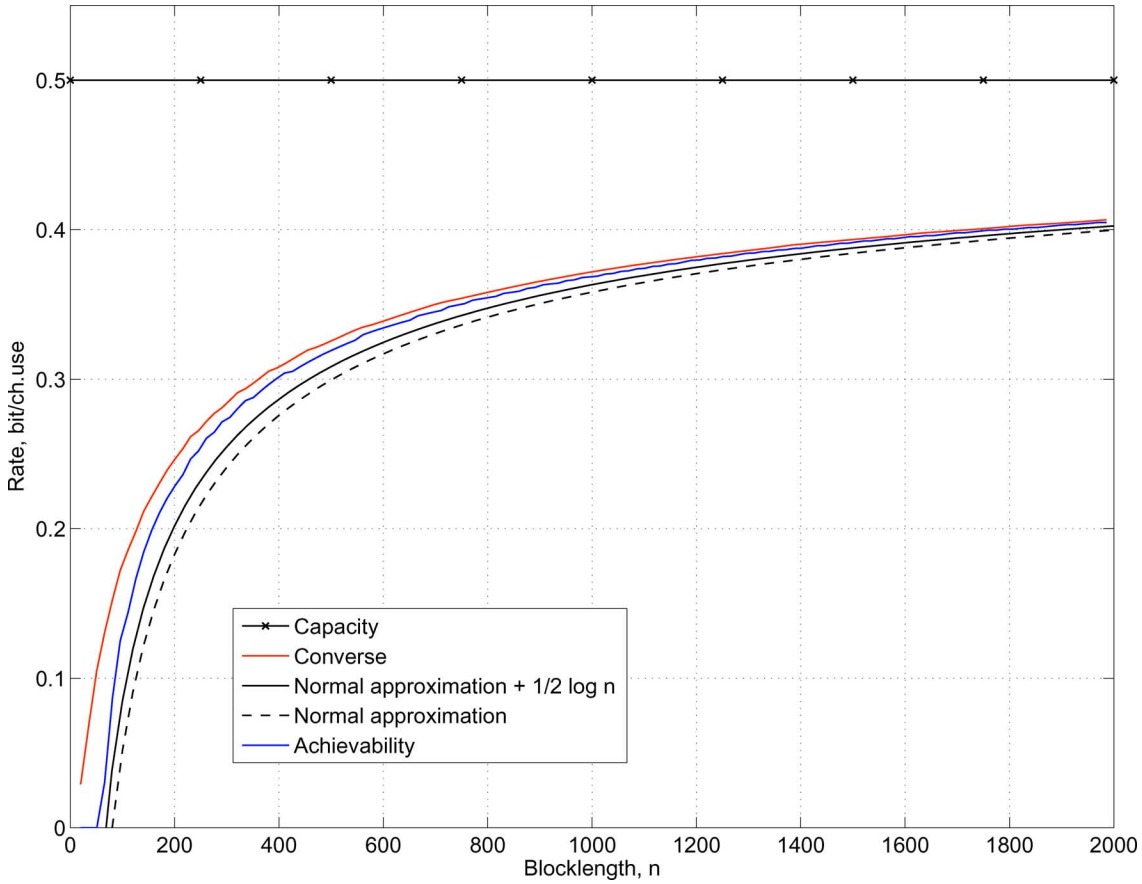


Fig. 9. Rate-blocklength tradeoff for the BSC with crossover probability  $\delta = 0.11$  and maximal block error rate  $\epsilon = 10^{-6}$ .

An interesting figure of merit for the AWGN channel is the excess energy per bit,  $\Delta E_b(n)$ , over that predicted by channel capacity incurred as a function of blocklength for a given required bit rate and block error rate:

$$\Delta E_b(n) = 10 \log_{10} \frac{P(n, R, \epsilon)}{\exp(2R) - 1} \quad (297)$$

where, according to the normal approximation,  $P(n, R, \epsilon)$  is the solution to

$$C - \sqrt{\frac{V}{n}} Q^{-1}(\epsilon) + \frac{1}{2n} \log n = R \quad (298)$$

and  $C$  and  $V$  are as in Theorem 54.

Fig. 14 gives a representative computation of (297)–(298) along with the corresponding lower<sup>18</sup> and upper bounds obtained from (208) and (220) respectively. We note a good precision of the simple approximation (297), e.g., for  $k = 100$  bits the gap to the achievability bound is only 0.04 dB. A similar comparison (without the normal approximation, of course) for rate 2/3 is presented in [48, Fig. 8].

*D. Application: Performance of Practical Codes*

How does the state-of-the-art compare against the finite blocklength fundamental limits? One such comparison is given in Fig. 12 where the lower curve depicts the performance of a certain family of multiedge low-density parity-check

(ME-LDPC) codes decoded via a low-complexity belief-propagation decoder [49]. We notice that in the absence of the nonasymptotic finite-blocklength curves, one has to compare the performance against the capacity alone. Such comparison leads to an incorrect conclusion that a given family of codes becomes closer to optimal with increasing blocklength. In reality we see that the relative gap to the finite blocklength fundamental limit is approximately constant. In other words, the fraction  $\frac{\log M_{LDPC}(n, \epsilon, P)}{\log M^*(n, \epsilon, P)}$  seems to be largely blocklength independent.

This observation leads us to a natural way of comparing two different codes over a given channel. Over the AWGN channel the codes have traditionally been compared in terms of  $E_b/N_0$ . Such comparison, although justified for a low-rate codes, unfairly penalizes higher rate codes. Instead, we define a normalized rate of a code with  $M$  codewords as (this can be extended to discrete channels parametrized by a scalar in a natural way)

$$R_{\text{norm}}(\epsilon) = \frac{\log M}{\log M^*(n, \epsilon, \gamma_{\min}(\epsilon))} \quad (299)$$

where  $\gamma_{\min}(\epsilon)$  is the smallest SNR at which the code still admits decoding with probability of error below  $\epsilon$ . The value  $\log M^*(n, \epsilon, \gamma_{\min}(\epsilon))$  can be safely replaced by an approximate value (296) with virtually no loss of precision for blocklength as low as 100.

The evolution of the coding schemes from 1980s (Voyager) to 2009 in terms of the normalized rate  $R_{\text{norm}}(10^{-4})$  is presented

<sup>18</sup>Another lower bound is given in [5, Fig. 3] which shows [3, (15)].

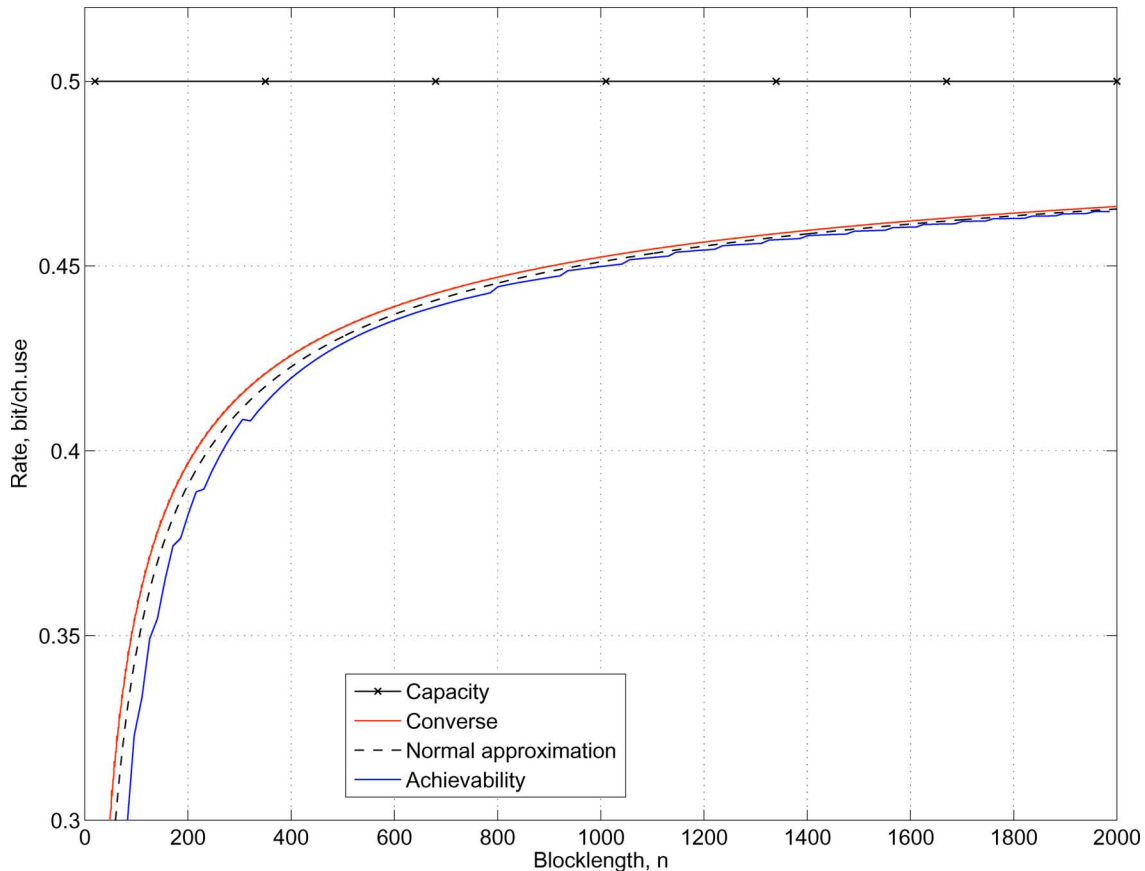


Fig. 10. Rate-blocklength tradeoff for the BEC with erasure probability  $\delta = 0.5$  and maximal block error rate  $\epsilon = 10^{-3}$ .

on Fig. 15. ME-LDPC is the same family as in [49, Fig. 12] and the rest of the data is taken from [5]. A comparison of certain turbo codes to Feinstein's bound and Shannon's converse can be found on [48, Figs. 6 and 7].

#### E. Application: Maximization of ARQ Throughput

A good analytical approximation to the maximal rate achievable with a given blocklength and error probability opens a variety of practical applications. In this subsection we consider a basic ARQ transmission scheme in which a packet is retransmitted until the receiver acknowledges successful decoding (which the receiver determines using a variety of known highly reliable hashing methods). Typically, the size  $k$  of the information packets is determined by the particular application, and both the blocklength  $n$  and the block error probability  $\epsilon$  are degrees of freedom. A natural objective is to maximize the average throughput (or, equivalently, minimize the average delivery delay) given by

$$T(k) = \max_{n, \epsilon} \frac{k}{n} (1 - \epsilon) \quad (300)$$

assuming decoding errors are independent for different retransmissions. The maximization in (300) is over those  $(n, \epsilon)$  such that

$$\log_2 M^*(n, \epsilon) = k. \quad (301)$$

Note that the number of required retransmissions is geometrically distributed, with mean equal to  $\frac{k}{T(k)}$ . In view of the tightness of the approximation in (223), it is sensible to maximize

$$\tilde{T}(k) = \max_n \frac{k}{n} \left[ 1 - Q \left( \frac{nC - k}{\sqrt{nV}} \right) \right] \quad (302)$$

where  $C$  and  $V$  are the channel capacity and channel dispersion, respectively. For the AWGN channel with SNR = 0 dB we show the results of the optimization in (302) in Fig. 16, where the optimal block error rate,  $\epsilon^*(k)$  is shown, and Fig. 17, where the optimal coding rate  $\frac{k}{n^*(k)}$  is shown. Table II shows the results of the optimization for the channel examples we have used throughout the paper. Of particular note is that for 1000 information bits, and a capacity-1/2 BSC, the optimal block error rate is as high as 0.0167.

The tight approximation to the optimal error probability as a function of  $k$  in Fig. 16 is the function

$$\tilde{\epsilon}(k) = \left( \frac{kC}{V} \ln \frac{kC}{2\pi V} \right)^{-1/2} \left( 1 - \frac{1}{\ln \frac{kC}{2\pi V}} \right) \quad (303)$$

obtained by retaining only the dominant terms in the asymptotic solution as  $k \rightarrow \infty$ .

## V. SUMMARY OF RESULTS

The main new nonasymptotic results shown in this paper are the following.

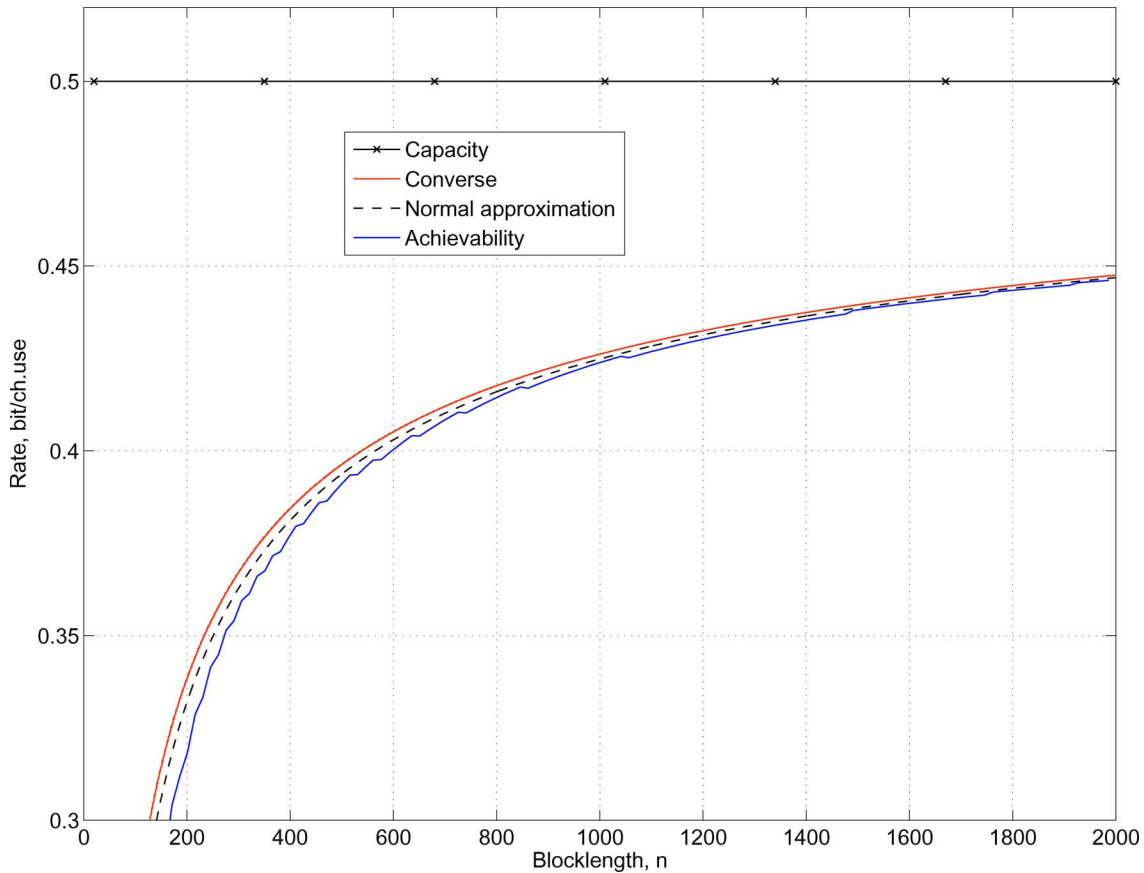


Fig. 11. Rate-blocklength tradeoff for the BEC with erasure probability  $\delta = 0.5$  and maximal block error rate  $\epsilon = 10^{-6}$ .

- 1) An exact expression (Theorem 15) for the error probability averaged over random codes which applies in full generality. In particular, it does not put any restrictions on the dependence of symbols within a codeword.
- 2) An upper bound (the RCU bound, Theorem 16) on the achievable average error probability for randomly chosen codes of a given size, which involves no bounding beyond the simple union bound. Loosening of the bound leads to the Shannon and Gallager bounds. When applied to a random ensemble, Poltyrev's BSC linear-code bound reduces to the RCU bound.
- 3) A simpler easier-to-compute bound (the DT bound, Theorem 17), which unlike previous achievability bounds contains no parameters to be optimized beyond the input distribution. The DT bound is tighter than the Shannon and Feinstein bounds, and, unlike the Gallager bound, it can be used to obtain the  $\sqrt{n}$  term in the normal approximation. For the BEC, the DT bound is generally tighter than the RCU bound. For channels with memory, in general the DT bound is easier to work with than any other new bounds in this paper; see [50].
- 4) A maximal error probability counterpart (Theorem 22) to the DT bound, obtained using the technique of sequential random coding.
- 5) The  $\kappa_i\beta$  bound (Theorem 25) which is a maximal error probability achievability bound based on the Neyman–Pearson lemma that uses an auxiliary output distribution. The  $\kappa_i\beta$  bound is particularly useful in the

setting of analog channels with cost constraints, and plays a key role in the normal approximation for the AWGN channel.

- 6) An auxiliary result (Theorem 26) which leads to a number of converse results, the most general of which is Theorem 27 which includes as simple corollaries the Fano inequality, the Wolfowitz converse and the Verdú–Han converse. Another corollary is Theorem 31 which can be viewed as a distillation of the essentials of the sphere-packing converse.
- 7) A tighter easy-to-compute converse bound (Theorem 38) for the BEC that holds even with noncausal feedback.

The tightness of the achievability bounds obtained by random coding is evidence that random-like codes (such as those arising in modern coding theory) not only achieve capacity but also do not sacrifice much performance for all but very short blocklengths. Numerical results with state-of-the-art codes show that about one half of the gap to capacity is due to the fundamental backoff due to finite blocklength; the other half of the gap is bridgeable with future advances in coding theory.

We have further shown the normal approximation to the maximal rate in the blocklength regime up to a term of  $O(\frac{\log n}{n})$  for both general discrete memoryless channels and additive white Gaussian noise channels, and up to  $O(\frac{1}{n})$  for both the BSC and the BEC. While for DMCs, the approach is a refinement of Strassen's [31], the Gaussian channel requires a different approach. The tightness of the approximation has been illustrated by comparison to the fixed-length bounds in Section III. It moti-

vates the use of the channel dispersion  $V$  (variance of the information density achieved by a capacity-achieving distribution), in conjunction with the channel capacity  $C$ , as a powerful analysis and design tool. In order to achieve a given fraction of capacity with a given error probability, the required blocklength is proportional to  $V/C^2$ .

The large deviations approach (reliability function) and the central-limit-theorem approach (dispersion) give a more refined analysis than that using only channel capacity. We note the following relationships and contrasts between both approaches:

- For rates near capacity, the reliability function behaves parabolically as

$$E(R) \approx \frac{(C - R)^2}{2V} \quad (304)$$

a fact that was known to Shannon as the unpublished, undated, unfinished manuscript [51] reproduced in Fig. 18 shows. Therefore, channel dispersion can be obtained by taking the second derivative of the reliability function at capacity. Since the reliability function is quite cumbersome to obtain for most channels, channel dispersion is far easier to obtain directly.

- According to the reliability function approximation, the blocklength required to sustain rate  $R = \eta C$  is inversely proportional to the reliability function evaluated at  $R$ , while according to the normal approximation it is proportional to

$$\frac{V}{C^2} \frac{1}{(1 - \eta)^2}.$$

Unless  $\eta$  is very close to 1 (in which case the factors are similar because of (304)) the normal approximation is substantially more accurate. In fact, even for rates substantially below capacity, the normal approximation remains accurate.

- Inherently, the large deviations approach does not capture the subexponential behavior (i.e., the “constant” factor in front of the exponential), which, for error probabilities and rates of practical interest, is more relevant than the inaccuracy of the approximation in (304).
- The reliability function approach predicts that the blocklength required to sustain  $\epsilon$  and a given desired rate scales with  $\log \frac{1}{\epsilon}$ , while the dispersion approach predicts that the scaling is  $(Q^{-1}(\epsilon))^2$ , which is equivalent for small  $\epsilon$  and rather more accurate otherwise.
- Often, the regime of very low  $\epsilon$  (the natural habitat of the reliability function), is not the desired one. Indeed, in many applications the error correcting code does not carry the whole burden of providing reliable communication; instead a protocol (such as ARQ) bootstraps a moderately low block error probability into very reliable communication (see Table II).
- For very low  $\epsilon$  neither approximation is accurate unless the blocklength is so high that the backoff from capacity is miniscule.

## APPENDIX A BOUNDS VIA LINEAR CODES

The goal of this appendix is to illustrate how Theorems 16 and 17, which give an upper bound on average probability of error, can also be used to derive an upper bound on maximal probability of error. To that end, we first notice that in both proofs we relied only on pairwise independence between randomly chosen codewords. So, the average probability of error for any other ensemble of codebooks with this property and whose marginals are identical and equal to  $P_X$  will still satisfy bounds of Theorems 16 and 17. In particular, for the BSC and the BEC we can generate an ensemble with equiprobable  $P_X$  by using a linear code with entries in its generating matrix chosen equiprobably on  $\{0, 1\}$ . Then, Theorems 16 and 17 guarantee the existence of the codebook, whose probability of error under ML decoding is small. Note that this is only possible if  $M = 2^k$  for some integer  $k$ . A question arises: for these structured codebooks are there randomized ML decoders whose maximal probability of error coincides with the average? This question is answered by the following result.

*Theorem 55:* Suppose that  $A$  is a group and suppose that there is a collection of measurable mappings  $T_x : B \mapsto B$  for each  $x \in A$  such that

$$P_{Y|X=x' \circ x} = P_{Y|X=x'} \circ (T_x)^{-1} \quad \forall x' \in A. \quad (305)$$

Then any code  $C$  that is a subgroup of  $A$  has a maximum likelihood decoder whose maximal probability of error coincides with the average probability of error.

Note that (305) can be reformulated as

$$\mathbb{E}[g(Y) | X = x' \circ x] = \mathbb{E}[g(T_x(Y)) | X = x'] \quad (306)$$

for all bounded measurable  $g : B \mapsto B$  and all  $x' \in A$ .

*Proof:* Define  $P_X$  to be a measure induced by the codebook  $C$

$$P_X(E) = \frac{1}{M} |E \cap C|. \quad (307)$$

Note that in this case  $P_Y$  induced by this  $P_X$  dominates all of  $P_{Y|X=x}$  for  $x \in C$

$$P_{Y|X=x} \ll P_Y \quad \forall x \in C. \quad (308)$$

Thus, we can introduce densities

$$f_{Y|X}(y|x) \triangleq \frac{dP_{Y|X=x}}{dP_Y}. \quad (309)$$

Observe that for any bounded measurable  $g$  we have

$$\mathbb{E}[g(Y)] = \mathbb{E}[g(T_x(Y))] \quad \forall x \in C. \quad (310)$$

Indeed

$$\mathbb{E}[g(T_x(Y))] = \sum_{x' \in C} \frac{1}{M} \mathbb{E}[g(T_x(Y)) | X = x'] \quad (311)$$

$$= \sum_{x' \in C} \frac{1}{M} \mathbb{E}[g(Y) | X = x' \circ x] \quad (312)$$

$$= \mathbb{E}[g(Y)] \quad (313)$$

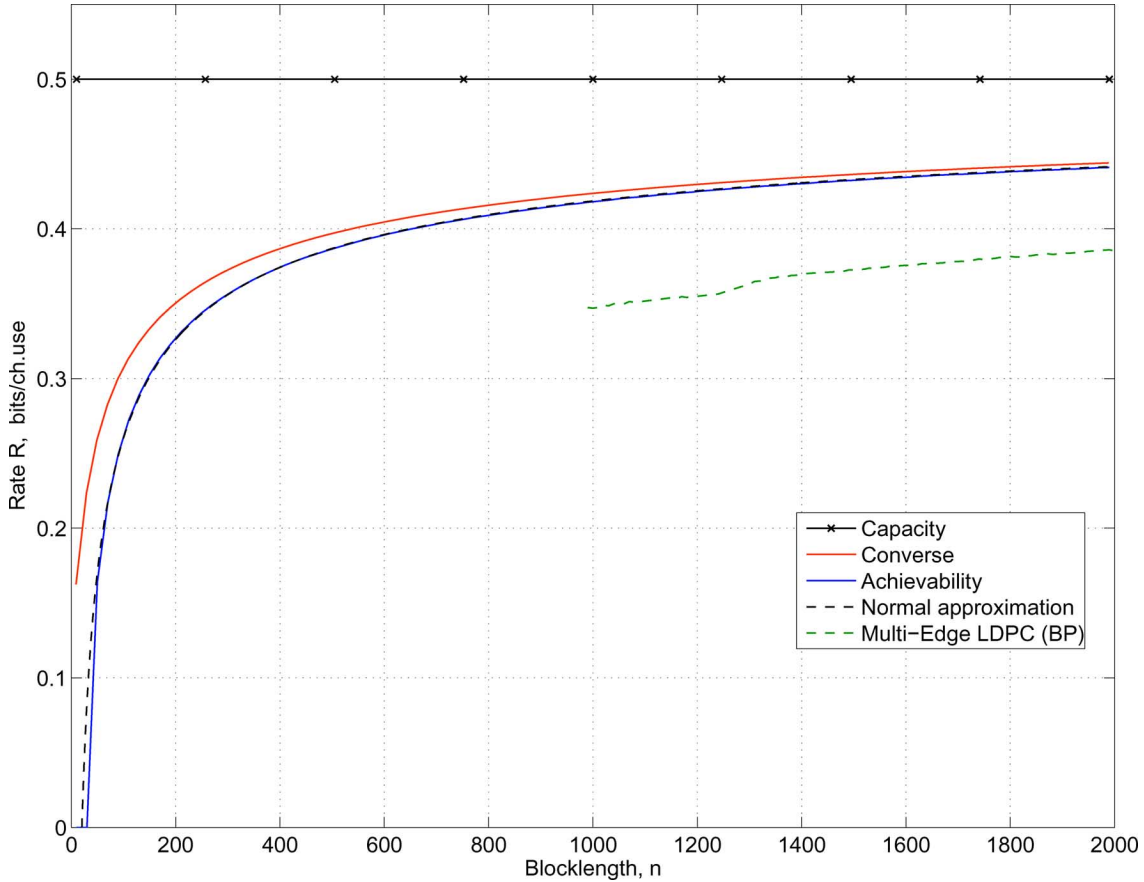


Fig. 12. Normal approximation for the AWGN channel,  $SNR = 0$  dB,  $\epsilon = 10^{-3}$ . The LDPC curve demonstrates the performance achieved by a particular family of multi-edge LDPC codes (designed by T. Richardson).

where (312) follows from (306). Also for any  $x, x' \in \mathcal{C}$  we have

$$f_{Y|X}(y|x') = f_{Y|X}(T_x(y)|x' \circ x) \quad P_Y\text{-a.s.} \quad (314)$$

Indeed, denote

$$E_1 = \{y : f_{Y|X}(y|x') < f_{Y|X}(T_x(y)|x' \circ x)\} \quad (315)$$

and assume that  $P_Y(E_1) = P_Y(T_x^{-1}E_1) > 0$ . Then, on one hand

$$P_{Y|X}[T_x^{-1}E_1 | x'] = \int_{\mathcal{B}} P_Y(dy) 1_{\{T_x(y) \in E_1\}} f_{Y|X}(y|x') \quad (316)$$

$$< \int P_Y(dy) 1_{\{T_x(y) \in E_1\}} f_{Y|X}(T_x(y)|x' \circ x) \quad (317)$$

$$= \int P_Y(dy) 1_{\{y \in E_1\}} f_{Y|X}(y|x' \circ x) \quad (318)$$

$$= P_{Y|X}[E_1 | x' \circ x] \quad (319)$$

where (318) follows from (310). But (319) contradicts (305) and hence  $P_Y(E_1) = 0$  and (314) is proved.

We proceed to define a decoder by the following rule: upon reception of  $y$  compute  $f_{Y|X}(y|x)$  for each  $x \in \mathcal{C}$ ; choose equiprobably among all the codewords that achieve the maximal  $f_{Y|X}(y|x)$ . Obviously, such decoder is maximum likelihood. We now analyze the conditional probability of error given

that the true codeword is  $x$ . Define two collections of functions of  $y$ , parameterized by  $x \in \mathcal{C}$

$$A_x(y) = \min \left\{ 1, \sum_{x' \in \mathcal{C}} 1_{\{f_{Y|X}(y|x') > f_{Y|X}(y|x)\}} \right\} \quad (320)$$

$$N_x(y) = \sum_{x' \in \mathcal{C}} 1_{\{f_{Y|X}(y|x') = f_{Y|X}(y|x)\}}. \quad (321)$$

It is easy to see that

$$\epsilon_x \triangleq \mathbb{P}[\text{error} | X = x] \quad (322)$$

$$= \mathbb{E} \left[ A_x(Y) + 1_{\{A_x(Y) = 0\}} \frac{N_x(Y) - 1}{N_x(Y)} \middle| X = x \right]. \quad (323)$$

If we denote the unit element of  $\mathcal{X}$  by  $x_0$ , then by (314) it is clear that

$$A_x \circ T_x = A_{x_0} \quad (324)$$

$$N_x \circ T_x = N_{x_0}. \quad (325)$$

But then, by (323) we have (326)–(329), shown at the bottom of the next page, where (326) follows because  $x_0$  is a unit of  $\mathcal{A}$ , (327) is by (306), and (328) is by (324) and (325). ■

The construction of  $T_x$  required in Theorem 55 is feasible for a large class of channels. For example, for an  $L$ -ary phase-shift-



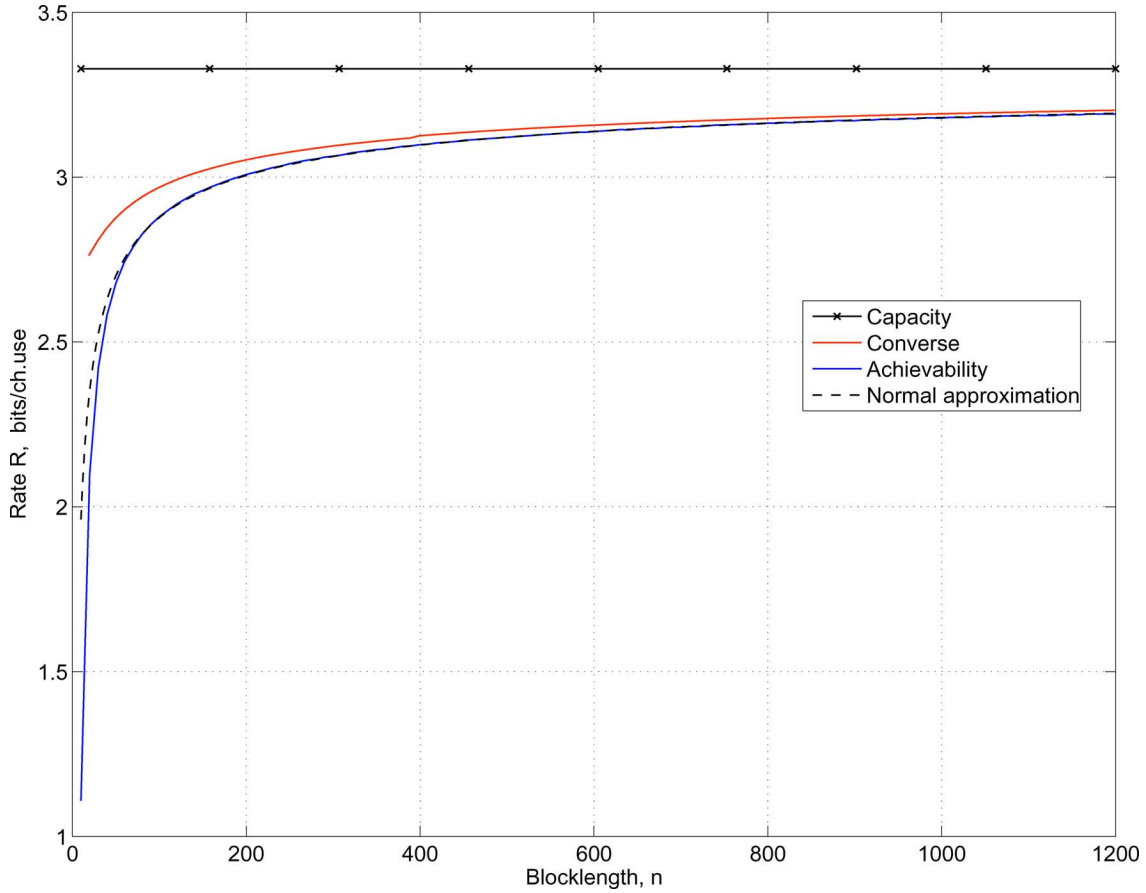


Fig. 13. Normal approximation for the AWGN channel,  $SNR = 20$  dB,  $\epsilon = 10^{-6}$ .

keying (PSK) modulated complex AWGN channel with soft decisions, we can assume that the input alphabet is  $\{e^{j\frac{2\pi k}{L}} k = 0, L - 1\}$ ; then

$$T_x(y) = yx \tag{330}$$

satisfies the requirements because  $P_{Y|X}(y|x')$  depends only on  $|y - x'|$  and  $|yx - x'x| = |y - x'|$ .

We give a general result for constructing  $T_x$ .

*Theorem 56:* Suppose that  $B$  is a monoid,  $A \subset B$  is a group (in particular  $A$  consists of only invertible elements of  $B$ ) and the channel is

$$Y = N \circ X \tag{331}$$

with  $N \in B$  being independent of  $X \in A$ . If each  $T_x(y) = y \circ x$  is measurable, then this family satisfies the conditions of Theorem 55.

*Proof:* Indeed, for any  $E \subset B$  we have

$$T_x^{-1}E = E \circ x^{-1}. \tag{332}$$

Then, on the one hand

$$P_{Y|X=x'}[E] = P_N[E \circ (x' \circ x)^{-1}] \tag{333}$$

but on the other hand

$$P_{Y|X=x'}[T_x^{-1}E] = P_{Y|X=x'}[E \circ x^{-1}] \tag{334}$$

$$= P_N[E \circ x^{-1} \circ x'^{-1}]. \tag{335}$$

■

$$\epsilon_x = \mathbb{E} \left[ A_x(Y) + 1\{A_x(Y) = 0\} \frac{N_x(Y) - 1}{N_x(Y)} \mid X = x_0 \circ x \right] \tag{326}$$

$$= \mathbb{E} \left[ A_x(T_x(Y)) + 1\{A_x(T_x(Y)) = 0\} \frac{N_x(T_x(Y)) - 1}{N_x(T_x(Y))} \mid X = x_0 \right] \tag{327}$$

$$= \mathbb{E} \left[ A_{x_0}(Y) + 1\{A_{x_0}(Y) = 0\} \frac{N_{x_0}(Y) - 1}{N_{x_0}(Y)} \mid X = x_0 \right] \tag{328}$$

$$= \epsilon_{x_0} \tag{329}$$

TABLE I  
BOUNDS ON THE MINIMAL BLOCKLENGTH  $n$  NEEDED TO ACHIEVE  $R = 0.9C$ .

| Channel                                | Converse      | RCU           | DT or $\kappa\beta$ | Error-exponent   | Normal Approx.   |
|--|---------------|---------------|---------------------|------------------|------------------|
| BEC(0.5), $\epsilon = 10^{-3}$         | $n \geq 899$  | $n \leq 1021$ | $n \leq 991$        | $n \approx 1380$ | $n \approx 955$  |
| BSC(0.11), $\epsilon = 10^{-3}$        | $n \geq 2985$ | $n \leq 3106$ | $n \leq 3548$       | $n \approx 4730$ | $n \approx 3150$ |
| AWGN, SNR = 0dB, $\epsilon = 10^{-3}$  | $n \geq 2550$ | $n \leq 2814$ | $n \leq 3400$       | $n \approx 4120$ | $n \approx 2750$ |
| AWGN, SNR = 20dB, $\epsilon = 10^{-6}$ | $n \geq 147$  | $n \leq 188$  | $n \leq 296$        | $n \approx 220$  | $n \approx 190$  |

It is easy to see that if we take  $\mathcal{A} = \mathbb{Z}_2$  and  $\mathbf{A} = \mathcal{A}^n$  then the BSC (even if the noise has memory) satisfies the conditions of Theorem 56. For the BEC we take  $\mathcal{A} = \{-1, 1\}$  and  $\mathcal{B} = \{-1, 0, 1\}$ , and the usual multiplication of reals converts  $\mathcal{B}$  to a monoid; taking the usual product  $\mathbf{A} = \mathcal{A}^n$  and  $\mathbf{B} = \mathcal{B}^n$  – we see that the BEC (even with memory) also satisfies the conditions of Theorem 56. Similar generalizations are possible for any additive noise channel with erasures.

#### APPENDIX B

##### NEYMAN–PEARSON LEMMA

*Lemma 57:* (For example, see [42]). Consider a space  $\mathcal{W}$  and probability measures  $P$  and  $Q$ . Then for any  $\alpha \in [0, 1]$  there exist  $\gamma > 0$  and  $\tau \in [0, 1)$  such that

$$\beta_\alpha(P, Q) = Q[Z_\alpha^* = 1] \quad (336)$$

and where<sup>19</sup> the conditional probability  $P_{Z^*|W}$  is defined via

$$Z_\alpha^*(W) = 1 \left\{ \frac{dP}{dQ} > \gamma \right\} + Z_\tau 1 \left\{ \frac{dP}{dQ} = \gamma \right\} \quad (337)$$

where  $Z_\tau \in \{0, 1\}$  equals 1 with probability  $\tau$  independent of  $W$ . The constants  $\gamma$  and  $\tau$  are uniquely determined by solving the equation

$$P[Z_\alpha^* = 1] = \alpha. \quad (338)$$

Moreover, any other test  $Z$  satisfying  $P[Z = 1] \geq \alpha$  either differs from  $Z_\alpha^*$  only on the set  $\left\{ \frac{dP}{dQ} = \gamma \right\}$  or is strictly larger with respect to  $Q$ :  $Q[Z = 1] > \beta_\alpha(P, Q)$ .

#### APPENDIX C

##### BINARY HYPOTHESIS TESTING: NORMAL APPROXIMATIONS

The next pair of results help us determine the asymptotic behavior of the optimal binary hypothesis tests with independent observations.

*Lemma 58:* Let  $\mathcal{A}$  be a measurable space with measures  $Q_i$  and  $P_i$ , with  $Q_i \ll P_i$  defined on it for  $i = 1, \dots, n$ . Define two measures on  $\mathcal{A}^n$ :  $Q = \prod_{i=1}^n Q_i$  and  $P = \prod_{i=1}^n P_i$ . Denote by  $\beta_\alpha$  the performance of the best randomized hypothesis test discriminating between  $Q$  and  $P$

$$\beta_\alpha = \inf_{P_{Z|Y^n}: Q[Z=1] \geq \alpha} \mathbb{P}[Z = 1]. \quad (339)$$

<sup>19</sup>In the case in which  $P$  is not absolutely continuous with respect to  $Q$ , we can define  $\frac{dP}{dQ}$  to be equal to  $+\infty$  on the singular set and hence to be automatically included in every optimal test.

Define

$$D_n = \frac{1}{n} \sum_{i=1}^n D(Q_i \| P_i) \quad (340)$$

$$\begin{aligned} V_n &= \frac{1}{n} \sum_{i=1}^n V(Q_i \| P_i) \\ &= \frac{1}{n} \sum_{i=1}^n \int \left( \log \frac{dQ_i}{dP_i} \right)^2 dQ_i - D(Q_i \| P_i)^2 \end{aligned} \quad (341)$$

$$T_n = \frac{1}{n} \sum_{i=1}^n \int \left| \log \frac{dQ_i}{dP_i} - D(Q_i \| P_i) \right|^3 dQ_i \quad (342)$$

$$B_n = 6 \frac{T_n}{V_n^{3/2}}. \quad (343)$$

Assume that all quantities are finite and  $V_n > 0$ . Then, for any  $\Delta > 0$

$$\begin{aligned} \log \beta_\alpha &\geq -nD_n - \sqrt{nV_n} Q^{-1} \left( \alpha - \frac{B_n + \Delta}{\sqrt{n}} \right) \\ &\quad + \log \Delta - \frac{1}{2} \log n \end{aligned} \quad (344)$$

$$\log \beta_\alpha \leq -nD_n - \sqrt{nV_n} Q^{-1} \left( \alpha + \frac{B_n}{\sqrt{n}} \right). \quad (345)$$

Each bound holds provided that the argument of  $Q^{-1}$  lies in  $(0, 1)$ .

*Proof of Lemma 58:* We will simply apply the Berry–Esseen Theorem 44 twice. We start from the lower bound. Observe that a logarithm of the Radon–Nikodym derivative  $\log \frac{dQ}{dP}$  is a sum of independent random variables by construction

$$\log \frac{dQ}{dP} = \sum_{i=1}^n \log \frac{dQ_i}{dP_i}. \quad (346)$$

Then applying (102), we have

$$\beta_\alpha \geq \frac{1}{\gamma_n} \left( \alpha - Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \right) \quad (347)$$

for  $\gamma_n > 0$ . Now set

$$\alpha_n = \alpha - \frac{B_n + \Delta}{\sqrt{n}} \quad (348)$$

which is positive since the argument of  $Q^{-1}$  in (344) is positive. Therefore, we let

$$\log \gamma_n = nD_n + \sqrt{nV_n} Q^{-1}(\alpha_n). \quad (349)$$

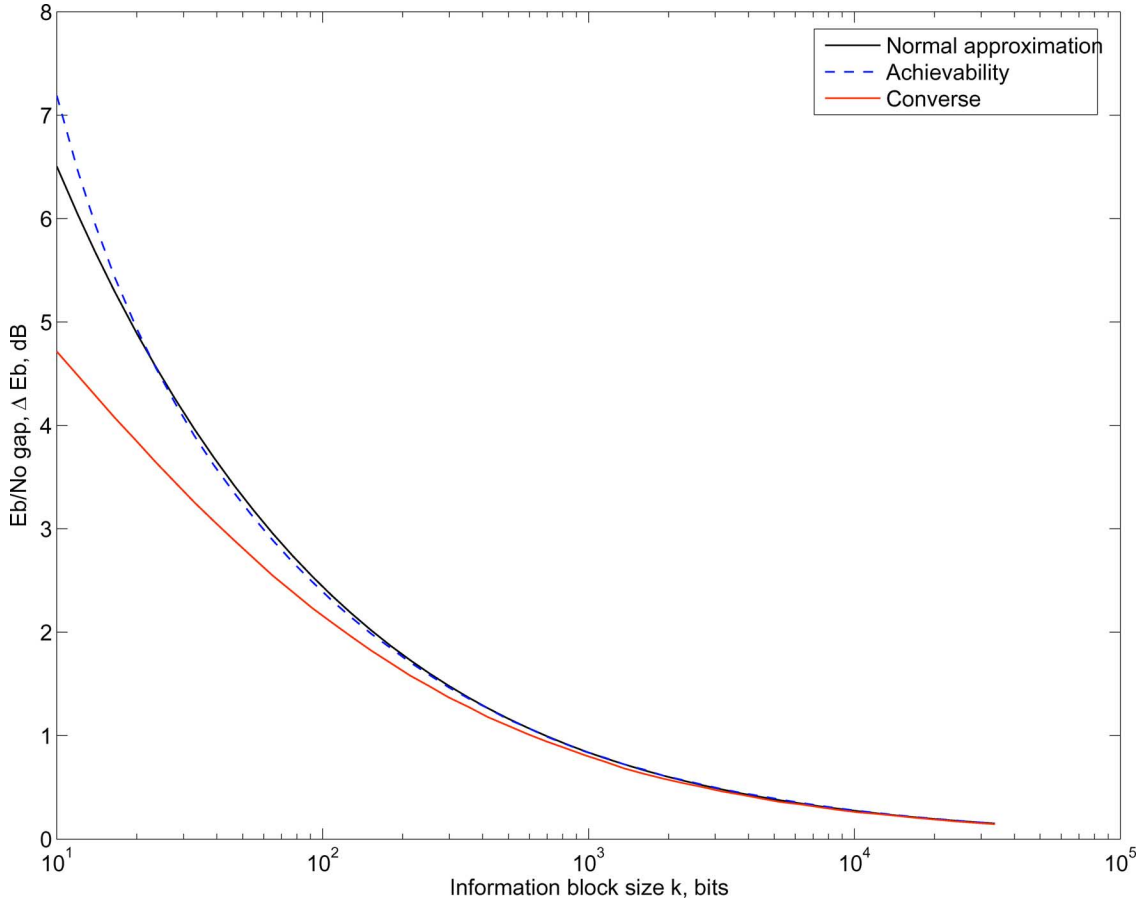


Fig. 14. Normal approximation for the  $\frac{E_b}{N_o}$  gap for the AWGN channel,  $R = 1/2$ ,  $\epsilon = 10^{-4}$ .

Then since  $\log \frac{dQ}{dP}$  is a sum of independent random variables, Theorem 44 applies and

$$\left| Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] - \alpha_n \right| \leq \frac{B_n}{\sqrt{n}}. \quad (350)$$

Consequently

$$Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \leq \alpha - \frac{\Delta}{\sqrt{n}}. \quad (351)$$

Substituting this bound into (347) we obtain (344).

For an upper bound, we use (103) which states that

$$\beta_\alpha^n \leq \frac{1}{\gamma_n} \quad (352)$$

whenever  $\gamma_n$  is such that

$$Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \geq \alpha. \quad (353)$$

Again, set

$$\alpha_n = \alpha + \frac{B_n}{\sqrt{n}} \quad (354)$$

which is strictly less than 1 since the argument of  $Q^{-1}$  in (345) is below 1. Similarly to (349) we choose

$$\log \gamma_n = nD + \sqrt{nV}Q^{-1}(\alpha_n). \quad (355)$$

From the Berry–Esseen bound, we have

$$\left| Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] - \alpha_n \right| \leq \frac{B_n}{\sqrt{n}}. \quad (356)$$

Consequently

$$Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \geq \alpha. \quad (357)$$

Thus, this choice of  $\gamma_n$  is valid for (352), and (345) follows. ■

Note that lower bound (344) holds only for  $n$  sufficiently large. A nonasymptotic bound is provided by the following result.

*Lemma 59:* In the notation of Lemma 58, we have

$$\log \beta_\alpha \geq -nD_n - \sqrt{\frac{2nV_n}{\alpha}} + \log \frac{\alpha}{2}. \quad (358)$$

*Proof:* Just as in the above argument, we start by writing

$$\beta_\alpha^n \geq \frac{1}{\gamma_n} \left( \alpha - Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \right). \quad (359)$$

We notice that

$$nD_n = \mathbb{E}_Q \left[ \log \frac{dQ}{dP} \right] \quad (360)$$

$$nV_n = \mathbb{E}_Q \left[ \left( \log \frac{dQ}{dP} - nD_n \right)^2 \right]. \quad (361)$$

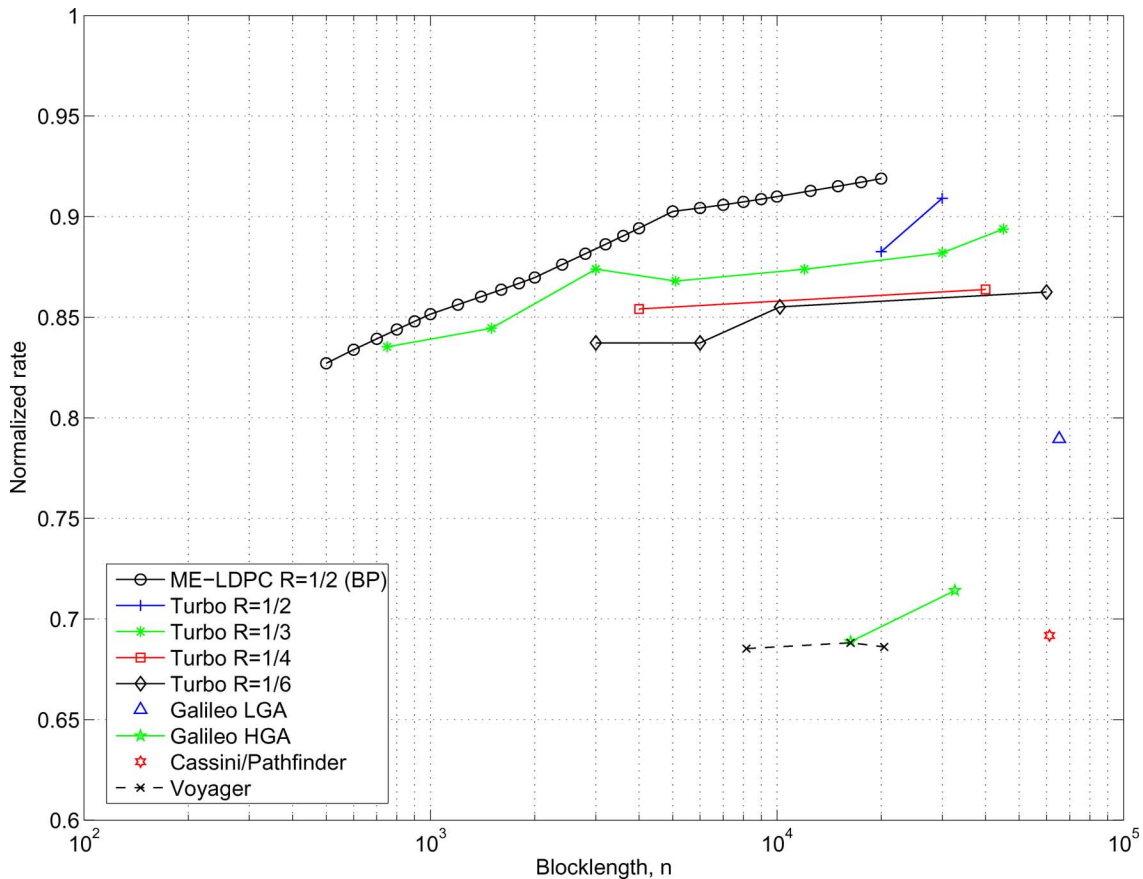


Fig. 15. Normalized rates for various practical codes over AWGN, probability of block error  $\epsilon = 10^{-4}$ .

Thus, if we set

$$\log \gamma_n = nD_n + \sqrt{\frac{2nV_n}{\alpha}} \quad (362)$$

then

$$Q \left[ \log \frac{dQ}{dP} \geq \log \gamma_n \right] \\ = Q \left[ \log \frac{dQ}{dP} - nD_n \geq \sqrt{\frac{2nV_n}{\alpha}} \right] \quad (363)$$

$$\leq Q \left[ \left( \log \frac{dQ}{dP} - nD_n \right)^2 \geq \frac{2nV_n}{\alpha} \right] \quad (364)$$

$$\leq \frac{\alpha}{2} \quad (365)$$

where (365) is by the Chebyshev inequality. Putting this into (359) we obtain the required result. ■

#### APPENDIX D EVALUATION OF $\kappa_n$ FOR THE AWGN CHANNEL

*Proof of Theorem 42:* According to Definition (107), we need to find the distribution  $P_{Z|Y^n}^*$  that, for every  $x \in \mathcal{F}_n$ , satisfies

$$\int_{\mathcal{B}} P_{Z|Y^n}^*(1|y) P_{Y^n|X^n=x}(dy) \geq \tau \quad (366)$$

and that attains the smallest possible value of

$$\int_{\mathcal{B}} P_{Z|Y^n}^*(1|y) P_{Y^n}(dy). \quad (367)$$

While, in general, this is a complex problem, the symmetry of the present case greatly simplifies the solution; we establish rigorously the spherical symmetry of the optimum attaining  $\kappa_n^n$ , and also suggest how to find symmetries in other (non-AWGN) problems of interest. We start by noting that any distribution  $P_{Z|Y^n}$  is completely determined by defining a function  $f : \mathcal{B} \mapsto [0, 1]$ , namely

$$f(y) = P_{Z|Y^n}(1|y). \quad (368)$$

We define the following class of functions on  $\mathcal{B} = \mathcal{B}^n$ :

$$\mathcal{F}_\tau = \left\{ f : \mathcal{B} \mapsto [0, 1] \right. \\ \left. \forall x \in \mathcal{A}^n : \int_{\mathcal{B}^n} f dP_{Y^n|X^n=x} \geq \tau \right\} \quad (369)$$

so that

$$\kappa^n(\tau) = \inf_{f \in \mathcal{F}_\tau} \int_{\mathcal{B}^n} f dP_{Y^n}. \quad (370)$$

Now we define another class, the subclass of spherically symmetric functions

$$\mathcal{F}_\tau^{\text{sym}} = \left\{ \phi \in \mathcal{F}_\tau : \phi(y) = \phi_r(\|y\|^2) \text{ for some } \phi_r \right\}. \quad (371)$$

We can then state the following.

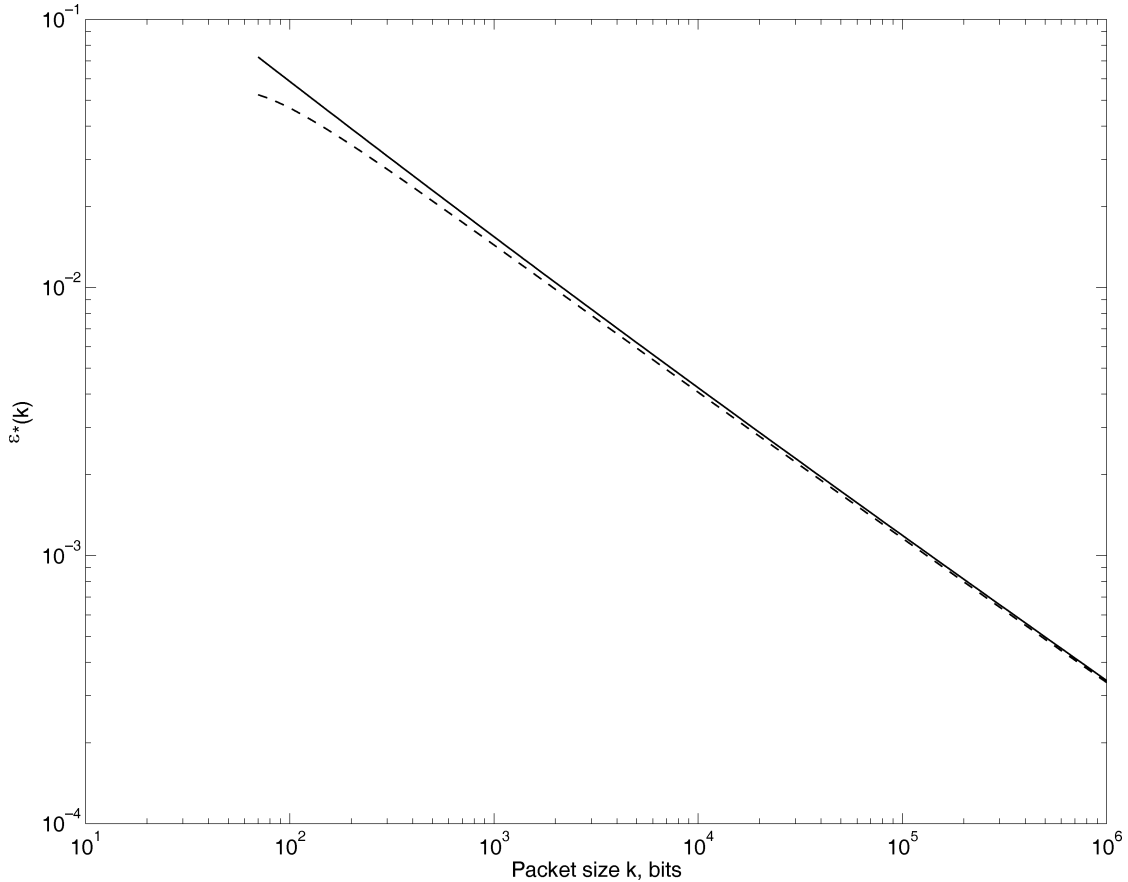


Fig. 16. Optimal block error rate  $\epsilon^*(k)$  maximizing average throughput under ARQ feedback for the AWGN channel with SNR = 0 dB. Solid curve is obtained by using normal approximation, dashed curve is an asymptotic formula (303).

*Lemma 60:* For the chosen  $P_{Y^n}$  and  $F_n$ , and for every  $\tau \in [0, 1]$  we have

$$\kappa_\tau^n(F_n, P_{Y^n}) = \inf_{\phi \in \mathcal{F}_\tau^{\text{sym}}} \int \phi dP_{Y^n}. \quad (372)$$

*Proof of Lemma 60:* The proof of Lemma 60 first defines a group  $G$  of transformations of  $\mathbb{B}$  (an orthogonal group  $O_n$  in this case) that permutes elements of the family of measures  $\{P_{Y^n}|_{X^n=x}, x \in F_n\}$  and that fixes  $P_Y^n$ . Then the optimum in the definition of  $\kappa_\tau^n$  can be sought as a function  $\mathbb{B} \mapsto [0, 1]$  that is constant on the orbits of  $G$  (this is the class  $\mathcal{F}_\tau^{\text{sym}}$ ).

Since  $\mathcal{F}_\tau^{\text{sym}} \subseteq \mathcal{F}_\tau$ , the inequality

$$\kappa_\tau^n \leq \inf_{\phi \in \mathcal{F}_\tau^{\text{sym}}} \int \phi dP_{Y^n} \quad (373)$$

is obvious. It remains to be shown that

$$\kappa_\tau^n \geq \inf_{\phi \in \mathcal{F}_\tau^{\text{sym}}} \int \phi dP_{Y^n}. \quad (374)$$

We will show that for every  $f \in \mathcal{F}_\tau$  there is a function  $\phi \in \mathcal{F}_\tau^{\text{sym}}$  with  $\int f dP_{Y^n} = \int \phi dP_{Y^n}$ . The claim (374) then follows trivially.

Define  $G$  to be the isometry group of a unit sphere  $\mathbb{S}^{n-1}$ . Then  $G = O(n)$ , the orthogonal group. Define a function on  $G \times G$  by

$$d(g, g') = \sup_{y \in \mathbb{S}^{n-1}} \|g(y) - g'(y)\|. \quad (375)$$

Since  $\mathbb{S}^{n-1}$  is compact,  $d(g, g')$  is finite. Moreover, it defines a distance on  $G$  and makes  $G$  a topological group. The group action  $H : G \times \mathbb{R}^n \mapsto \mathbb{R}^n$  defined as

$$H(g, y) = g(y) \quad (376)$$

is continuous in the product topology on  $G \times \mathbb{R}^n$ . Also,  $G$  is a separable metric space. Thus, as a topological space, it has a countable basis. Consequently, the Borel  $\sigma$ -algebra on  $G \times \mathbb{R}^n$  coincides with the product of Borel  $\sigma$ -algebras on  $G$  and  $\mathbb{R}^n$ :

$$\mathcal{B}(G \times \mathbb{R}^n) = \mathcal{B}(G) \times \mathcal{B}(\mathbb{R}^n). \quad (377)$$

Finally,  $H(g, y)$  is continuous and hence is measurable with respect to  $\mathcal{B}(G \times \mathbb{R}^n)$  and thus is also a measurable mapping with respect to a product  $\sigma$ -algebra.

It is also known that  $G$  is compact. On a compact topological group there exists a unique (right Haar) probability measure  $\mu$  compatible with the Borel  $\sigma$ -algebra  $\mathcal{B}(G)$ , and such that

$$\mu(Ag) = \mu(A) \quad \forall g \in G, A \in \mathcal{B}(G). \quad (378)$$

Now take any  $f \in \mathcal{F}_\tau$  and define an averaged function  $\phi(y)$  as

$$\phi(y) \triangleq \int_G (f \circ H)(g, y) \mu(dg). \quad (379)$$

Note that as shown above  $f \circ H$  is a positive measurable mapping  $G \times \mathbb{B} \mapsto \mathbb{R}_+$  with respect to corresponding Borel  $\sigma$ -algebras.

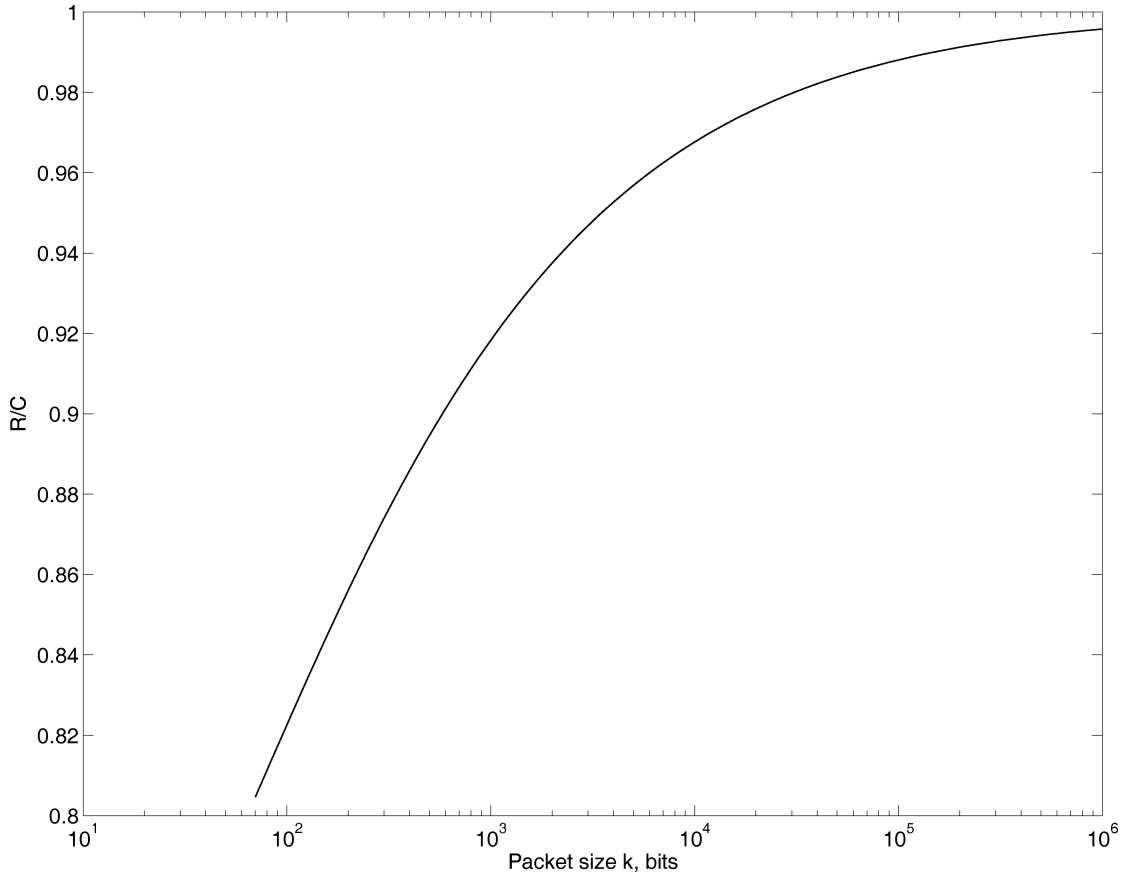


Fig. 17. Optimal coding rate maximizing average throughput under ARQ feedback for the AWGN channel with SNR = 0 dB. Solid curve is obtained using normal approximation.

Then by Fubini's theorem, the function  $\phi : \mathcal{B} \mapsto \mathbb{R}_+$  is also positive measurable. Moreover

$$0 \leq \phi(y) \leq \int_G 1 \mu(dg) = 1. \quad (380)$$

Define for convenience

$$Q_Y^x \triangleq P_{Y^n | X^n = x}. \quad (381)$$

Then

$$\int_{\mathcal{B}} \phi(y) Q_Y^x(dy) = \int_{\mathcal{B}} Q_Y^x(dy) \int_G (f \circ H)(g, y) \mu(dg) \quad (382)$$

$$= \int_G \mu(dg) \int_{\mathcal{B}} (f \circ H)(g, y) Q_Y^x(dy). \quad (383)$$

Change of the order of integration is possible by Fubini's theorem because  $f \circ H$  is a bounded function. By the change of variable formula

$$\begin{aligned} \int_G \mu(dg) \int_{\mathcal{B}} (f \circ g)(y) Q_Y^x(dy) \\ = \int_G \mu(dg) \int_{\mathcal{B}} f(Q_Y^x \circ g^{-1})(dy). \end{aligned} \quad (384)$$

By the definition of  $Q_Y^x$  we have, for every measurable set  $E$ ,  $Q_Y^x[E] = Q_Y^0[E - x]$  and the measure  $Q_Y^0$  is fixed under all isometries of  $\mathbb{R}^n$ :

$$\forall g \in G : Q_Y^0[F] = Q_Y^0[g(F)]. \quad (385)$$

But then

$$(Q_Y^x \circ g^{-1})[E] \triangleq Q_Y^x[g^{-1}(E)] \quad (386)$$

$$= Q_Y^0[g^{-1}(E) - x] \quad (387)$$

$$= Q_Y^0[g^{-1}[E - g(x)]] \quad (388)$$

$$= Q_Y^{g(x)}[E]. \quad (389)$$

This proves that

$$Q_Y^x \circ g^{-1} = Q_Y^{g(x)}. \quad (390)$$

It is important that  $x \in F_n$  implies  $g(x) \in F_n$ . In general terms, without AWGN channel specifics, the above argument shows that in the space of all measures on  $\mathcal{B}$  the subset  $\{Q_Y^x, x \in F_n\}$  is invariant under the action of  $G$ .

But  $f \in \mathcal{F}_\tau$  and thus  $\int f dQ_Y^x \geq \tau$  for every  $x \in F_n$ . So, from (384) and (390), we conclude

$$\int_{\mathcal{B}} \phi dQ_Y^x \geq \int_G \tau \mu(dg) = \tau. \quad (391)$$

TABLE II  
OPTIMAL BLOCK ERROR RATE FOR PACKET SIZE  $k = 1000$  BITS

| Channel          | Optimal $\epsilon^*(k)$ | Optimal $R/C$ | Optimal throughput |
|------------------|-------------------------|---------------|--------------------|
| BEC(0.5)         | $8.1 \cdot 10^{-3}$     | 0.95          | 0.94               |
| BSC(0.11)        | $16.7 \cdot 10^{-3}$    | 0.91          | 0.90               |
| AWGN, SNR = 0dB  | $15.5 \cdot 10^{-3}$    | 0.92          | 0.90               |
| AWGN, SNR = 20dB | $6.2 \cdot 10^{-3}$     | 0.96          | 0.95               |

Together with (380) this establishes that  $\phi \in \mathcal{F}_\tau$ . Now, the  $P_{Y^n}$  measure is also fixed under any  $g \in G$ :

$$P_{Y^n} \circ g^{-1} = P_{Y^n}. \quad (392)$$

Then replacing  $Q_Y^x$  with  $P_{Y^n}$  in (384), we obtain

$$\int_{\mathbb{B}} \phi dP_{Y^n} = \int_G \mu(dg) \int_{\mathbb{B}} f(y) (P_{Y^n} \circ g^{-1})(dy) \quad (393)$$

$$= \int_{\mathbb{B}} f dP_{Y^n}. \quad (394)$$

It remains to show that  $\phi \in \mathcal{F}_\tau^{\text{sym}}$ ; but, this is a simple consequence of the choice of  $\mu$ . Indeed for any  $g' \in G$

$$(\phi \circ g')(y) = \int_G (f \circ H)(g, g'(y)) \mu(dg) \quad (395)$$

$$= \int_G (f \circ H)(gg', y) \mu(dg) \quad (396)$$

$$= \int_G (f \circ H)(g'', y) \mu(dg'') \quad (397)$$

$$= \phi(y). \quad (398)$$

In the last equality we used a change of measure and the invariance of  $\mu$  under right translations. Thus,  $\phi$  must be constant on the orbits of  $G$  and hence, depends only on the norm of  $y$ . To summarize, we have shown that  $\phi$  belongs to  $\mathcal{F}_\tau^{\text{sym}}$  and

$$\int \phi dP_{Y^n} = \int f dP_{Y^n}. \quad (399)$$

■

*Proof of Theorem 42 (Continued):* By Lemma 60 we obtain a value of  $\kappa_\tau^n$  by optimizing over spherically symmetric functions. First, we will simplify the constraints on the functions in  $\mathcal{F}_\tau^{\text{sym}}$ . Define  $Q_Y^x$  and  $G$  as in the proof of Lemma 60. As we saw in that proof, each transformation  $g \in G$  carries one measure  $Q_Y^x$  into another  $Q_Y^{x'}$ . Also  $x' = g(x)$  in this particular case, but this is not important. What is important, however, is that if  $x \in F_n$  then  $x' \in F_n$ . If we define

$$\mathcal{Q} = \{Q_Y^x, x \in F_n\} \quad (400)$$

then, additionally, the action of  $G$  on  $\mathcal{Q}$  is transitive. This opens the possibility that the system of constraints on  $\phi \in \mathcal{F}_\tau^{\text{sym}}$  might be overdetermined. Indeed, suppose that  $\phi$  satisfies

$$\int_{\mathbb{B}} \phi dQ_0 \geq \tau \quad (401)$$

for some  $Q_0 \in \mathcal{Q}$ . Then for any measure  $Q \in \mathcal{Q}$  there is a transformation  $g \in G$  such that

$$Q = Q_0 \circ g^{-1}. \quad (402)$$

But then

$$\int_{\mathbb{B}} \phi dQ = \int_{\mathbb{B}} \phi \circ g dQ_0 = \int_{\mathbb{B}} \phi dQ_0. \quad (403)$$

Here the last equality follows from the fact that all members of  $\mathcal{F}_\tau^{\text{sym}}$  are spherically symmetric functions and as such are fixed under  $G$ :  $\phi \circ g = \phi$ . That is, once a symmetric  $\phi$  satisfies

$$\int_{\mathbb{B}} \phi dP_{Y^n|X^n=x_0} \geq \tau \quad (404)$$

for one  $x_0 \in F_n$ , it automatically satisfies the same inequality for all  $x \in F_n$ . So we are free to check (404) at one arbitrary  $x_0$  and then conclude that  $\phi \in \mathcal{F}_\tau^{\text{sym}}$ . For convenience we choose

$$x_0 = (\sqrt{P}, \sqrt{P}, \dots, \sqrt{P}). \quad (405)$$

Since all functions in  $\mathcal{F}_\tau^{\text{sym}}$  are spherically symmetric we will work with their radial parts:

$$\phi(y) = \phi_r(\|y\|^2). \quad (406)$$

Note that  $P_{Y^n}$  induces a certain distribution on  $R = \|Y^n\|^2$ , namely

$$P_0 \sim \sum_{i=1}^n (1+P)Z_i^2 \quad (407)$$

(as previously the  $Z_i$ 's denote i.i.d. standard Gaussian random variables). Similarly,  $P_{Y^n|X^n=x_0}$  induces a distribution on  $R = \|Y^n\|^2$ , namely,

$$P_1 \sim \sum_{i=1}^n (Z_i + \sqrt{P})^2. \quad (408)$$

Finally, we see that  $\kappa_\tau^n$  is

$$\kappa_\tau^n = \inf_{\{\phi_r: \int \phi_r dP_1 \geq \tau\}} \int \phi_r dP_0 \quad (409)$$

—a randomized binary hypothesis testing problem with  $P_1$ (decide  $P_1$ )  $\geq \tau$ .

Finally, we are left to note that the existence of a unique optimal solution  $\phi_r^*$  is guaranteed by the Neyman–Pearson lemma

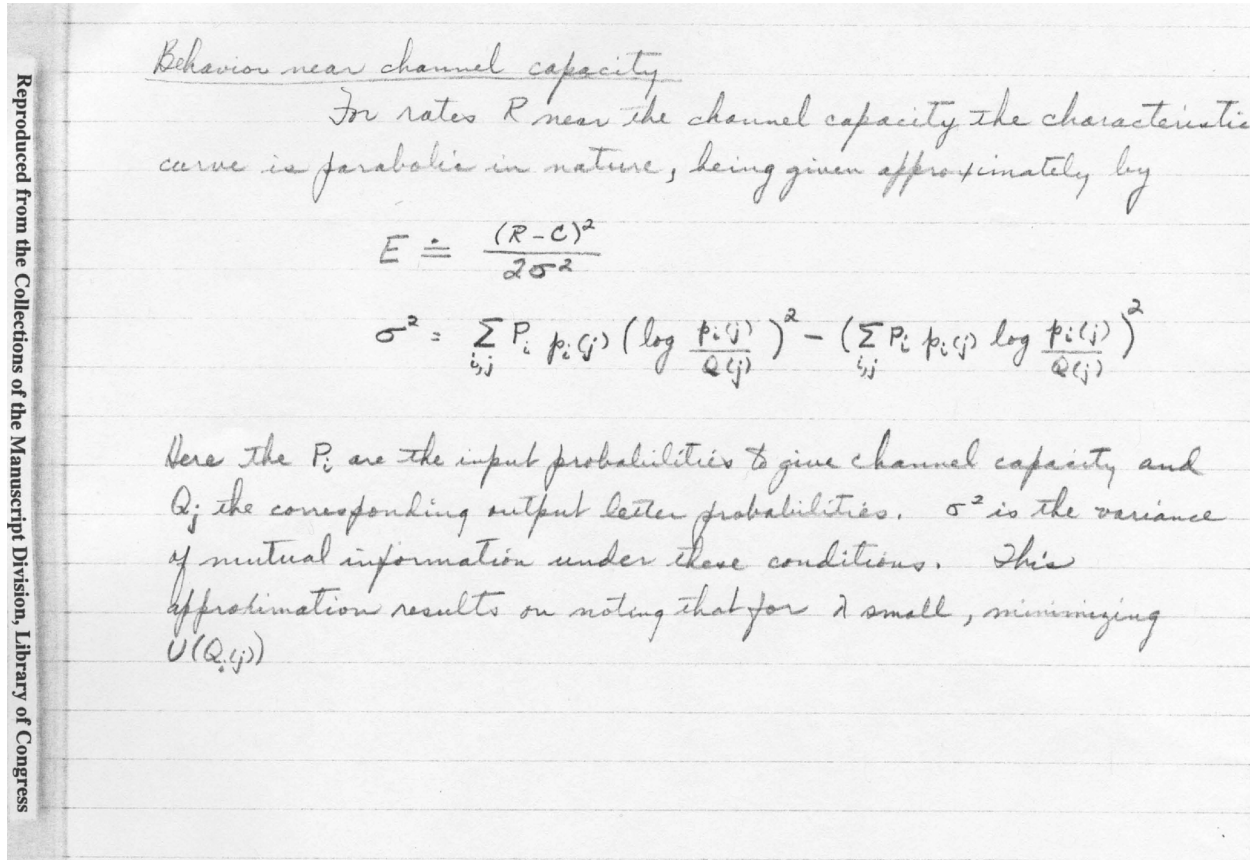


Fig. 18. [51] reproduced in its entirety.

(Appendix B). To conclude the proof we must show that the solution of (211) exists and thus that  $\phi_r^*$  is an indicator function (i.e., there is no “randomization on the boundary” of a likelihood ratio test). To that end, we need to show that for any  $\gamma$  the set

$$A_\gamma = \left\{ \frac{p_1(r)}{p_0(r)} = \gamma \right\} \quad (410)$$

satisfies  $P_1[A_\gamma] = 0$ . To show this, we will first show that each set  $\{A_\gamma \cap [0, K]\}$  is finite; then, its Lebesgue measure is zero, and since  $P_1$  is absolutely continuous with respect to Lebesgue measure we conclude from the monotone convergence theorem that

$$P_1[A_\gamma] = \lim_{K \rightarrow \infty} P_1[A_\gamma \cup [0, K]] = 0. \quad (411)$$

Note that the distribution  $P_0$  is a scaled  $\chi^2$ -distribution with  $n$  degrees of freedom; thus (e.g., [43, (26.4.1)]) the PDF of  $P_0$  is indeed given by (212). The distribution  $P_1$  is the noncentral  $\chi^2$ -distribution with  $n$  degrees of freedom and noncentrality parameter,  $\lambda$ , equal to  $nP$ . Thus (see [43, (26.4.25)]) we can write the PDF of  $P_1$  as expressed in (213). Using these expressions we obtain

$$f(r) \doteq \frac{p_1(r)}{p_0(r)} = e^{-\mu r} \sum_{i=0}^{\infty} a_i r^i. \quad (412)$$

The coefficients  $a_i$  are such that the series converges for any  $r < \infty$ . Thus, we can extend  $f(r)$  to be an analytic function  $F(z)$  over the entire complex plane. Now fix a  $K \in (0, \infty)$  and denote

$$S = A_\gamma \cap [0, K] = f^{-1}\{\gamma\} \cap [0, K]. \quad (413)$$

By the continuity of  $f$  the set  $S$  is closed. Thus,  $S$  is compact. Suppose that  $S$  is infinite; then there is sequence  $r_k \in S$  converging to some  $r^* \in S$ . But then from the uniqueness theorem of complex analysis, we conclude that  $F(z) = \gamma$  over the entire disk  $|z| \leq K$ . Since  $f(r)$  cannot be constant, we conclude that  $S$  is finite. ■

To enable non-AWGN applications of the  $\kappa\beta$  bound, let us summarize the general ideas used to prove Lemma 60 and Theorem 42. The proof of Lemma 60 first defines a group  $G$  of transformations of  $\mathcal{B}$  (an orthogonal group  $O_n$  in this case) that permutes elements of the family of measures  $\{P_{Y^n|X^n=x}, x \in \mathcal{F}_n\}$  and that fixes  $P_{Y^n}$ . Then the optimum in the definition of  $\kappa_\tau^n$  can be sought as a function  $\mathcal{B} \mapsto [0, 1]$  that is constant on the orbits of  $G$  (this is the class  $\mathcal{F}_\tau^{\text{sym}}$ ). Carrying this idea forward, in the proof of Theorem 42 we note that the action of  $G$  on  $\{P_{Y^n|X^n=x}, x \in \mathcal{F}_n\}$  is transitive and thus a set of conditions on  $\phi \in \mathcal{F}_\tau^{\text{sym}}$  can be replaced by just one

$$\int \phi dP_{Y^n|X^n=x_0} \geq \tau \quad (414)$$



for any  $x_0$ . If  $x_0$  is conveniently chosen, then computation of  $\kappa_\tau^n$  is a matter of solving a single randomized binary hypothesis testing problem between two memoryless distributions.

*Lemma 61:* For every  $P > 0$  there are constants  $C_1 > 0$  and  $C_2 > 0$  such that for all sufficiently large  $n$  and all  $\tau \in [0, 1]$ ,

$$\kappa_\tau^n \geq \frac{1}{C_1} (\tau - e^{-C_2 n}). \tag{415}$$

*Proof:* Recall that  $\kappa_\tau^n$  is determined by a binary hypothesis testing problem between  $P_0^{(n)}$  and  $P_1^{(n)}$ , as defined by (407) and (408). We will omit indices  $(n)$  where it does not cause confusion. Also in this proof all exp exponents are to the base  $e$ . The argument consists of two steps.

Step 1) There is a  $\delta > 0$  such that for all  $n \geq 1$  the Radon–Nikodym derivative  $\frac{dP_1^{(n)}}{dP_0^{(n)}}(r)$  is upper bounded by a constant  $C_1$  on the set

$$r \in R_n \triangleq [n(1 + P - \delta), n(1 + P + \delta)]. \tag{416}$$

Step 2) Since the measures  $P_1^{(n)}$  have mean  $n(1 + P)$ , by the Chernoff bound there is a constant  $C_2$  such that

$$P_1^{(n)}[\{R_n\}^c] \leq e^{-C_2 n}. \tag{417}$$

Now choose any set  $A$  such that  $P_1[A] \geq \tau$ . Then

$$P_1[A \cap R_n] \geq P_1[A] - P_1[\{R_n\}^c] \tag{418}$$

$$\geq \tau - e^{-C_2 n}. \tag{419}$$

But then

$$P_0[A] \geq P_0[A \cap R_n] \tag{420}$$

$$= \int_{\mathbb{R}_+} 1_{A \cap R_n} dP_0 \tag{421}$$

$$= \int_{A \cap R_n} \frac{dP_0}{dP_1} dP_1 \tag{422}$$

$$\geq \frac{1}{C_1} \int_{A \cap R_n} dP_1 \tag{423}$$

$$\geq \frac{1}{C_1} (\tau - e^{-C_2 n}). \tag{424}$$

This establishes the required inequality. The rest is devoted to proving Step 1, namely

$$f_n(r) \triangleq \frac{dP_1}{dP_0} \leq C_1 \quad \text{on } R_n \quad \forall n. \tag{425}$$

We have already discussed some properties of  $f_n(r)$  in (412). Here, however, we will need a precise expression for it, easily obtainable via (212) and (213):

$$f_n(r) = (1 + P)^{n/2} \exp\left\{-n\frac{P}{2} - r\frac{P}{2P+2}\right\} \times (nP r)^{-n/4+1/2} 2^{n/2} \Gamma\left(\frac{n}{2}\right) I_{n/2-1}(\sqrt{nPr}) \tag{426}$$

where  $I_{n/2-1}(x)$  is the modified Bessel function of the first kind.

We will consider only the case in which  $n$  is even. This is possible because in [44] it is shown that

$$\mu > \nu \geq 0 \implies I_\mu(x) < I_\nu(x) \tag{427}$$

for all  $x > 0$ . Thus, if  $n$  is odd then an upper bound is obtained by replacing  $I_{n/2-1}$  with  $I_{n/2-3/2}$ . Now for integer index  $k = n/2 - 1$  the following bound is shown in [45]:

$$I_k(z) \leq \sqrt{\frac{\pi}{8}} e^z \frac{1}{\sqrt{z}} \left(1 + \frac{k^2}{z^2}\right)^{-1/4} \times \exp\left\{-k \sinh^{-1} \frac{k}{z} + z \left(\sqrt{1 + \frac{k^2}{z^2}} - 1\right)\right\}. \tag{428}$$

Note that we need to establish the bound only for  $r$ 's that are of the same order as  $n$ ,  $r = O(n)$ . Thus, we will change the variable

$$r = nt \tag{429}$$

and seek an upper bound on  $f_n(nt)$  for all  $t$  inside some interval containing  $(1 + P)$ .

Using (428) and the expression

$$\ln \Gamma\left(\frac{n}{2}\right) = \frac{n-1}{2} \ln \frac{n}{2} - \frac{n}{2} + O(1) \tag{430}$$

$f_n(r)$  in (426) can be upper bounded, after some algebra, as

$$f_n(nt) \leq \exp\left\{-\frac{n}{2} K(t, P) + O(1)\right\}. \tag{431}$$

Here the  $O(1)$  term is uniform in  $t$  for all  $t$  on any finite interval not containing zero, and

$$K(t, P) = -\ln\left\{1 + \sqrt{1 + 4Pt}\right\} + \sqrt{1 + 4Pt} + \ln(1 + P) - P - \frac{Pt}{P+1} - 1 + \ln 2. \tag{432}$$

A straightforward exercise shows that a maximum of  $K(t, P)$  is attained at  $t^* = 1 + P$  and

$$K_{\max} = K(t^*, P) = 0. \tag{433}$$

Thus

$$f_n(nt) \leq O(1) \quad t \in [a, b] \quad a > 0. \tag{434}$$

In particular (425) holds if we take, for example,  $a = (1 + P) - 1$  and  $b = (1 + P) + 1$ . ■

In fact, the Radon–Nikodym derivative is bounded for all  $r$ , not only  $r \in R_n$  and, hence

$$\kappa_\tau^n \geq \frac{\tau}{C_1} \tag{435}$$

instead of the weaker (415). But showing that this holds for all  $r$  complicates the proof unnecessarily.

APPENDIX E  
CONTINUITY ON THE SIMPLEX

Our first results are concerned with properties of  $U(P, W)$  and  $V(P, W)$ .

*Lemma 62:* The functions  $U(P, W)$ ,  $V(P, W)$  and  $T(P, W)$  are continuous on  $\mathcal{P}$ . Functions  $U(P, W)$  and  $V(P, W)$  coincide on  $\Pi$ . The following inequality holds:

$$V(P, W) \leq U(P, W) \leq 2g(\min\{|\mathcal{A}|, |\mathcal{B}|\}) - [I(P, W)]^2 \quad (436)$$

where

$$g(n) = \begin{cases} 0.6 \log^2 e & n = 2, \\ \log^2 n & n \geq 3. \end{cases} \quad (437)$$

*Proof:* First, note that  $U(P, W)$ ,  $V(P, W)$  and  $T(P, W)$  are well-defined and finite. Indeed, each one is a sum of finitely many terms. We must show that every term is well-defined. This is true since, whenever  $W(y|x) = 0$  or  $PW(y) = 0$  or  $P(x) = 0$ , we have  $P(x)W(y|x) = 0$  and thus

$$P(x)W(y|x) \left[ \log \frac{W(y|x)}{PW(y)} \right]^2$$

and

$$P(x)W(y|x) \left| \log \frac{W(y|x)}{PW(y)} - D(W_x \| PW) \right|^3$$

are both equal to zero by convention. On the other hand, if  $P(x) > 0$  then  $W_x \ll PW$  and thus  $D(W_x \| PW)$  is a well-defined finite quantity. Second, take a sequence  $P_n \rightarrow P$ . Then we want to prove that each term in  $U(P, W)$  is continuous, i. e.,

$$\begin{aligned} P_n(x)W(y|x) \left[ \log \frac{W(y|x)}{P_n W(y)} \right]^2 \\ \rightarrow P(x)W(y|x) \left[ \log \frac{W(y|x)}{PW(y)} \right]^2. \end{aligned} \quad (438)$$

If  $W(y|x) = 0$  then this is obvious. If  $P_n(x) \neq 0$  then this is also true since the argument of the logarithm is bounded away from 0 and  $+\infty$ . So, we assume  $P_n(x) \rightarrow 0$  and we must show that then the complete quantity also tends to 0. For  $P_n(x) > 0$  we notice that

$$\log(P_n(x)W(y|x)) \leq \log P_n W(y) \leq 0. \quad (439)$$

Thus,

$$\begin{aligned} |\log W(y|x) - \log P_n W(y)|^2 \\ \leq 2(\log^2 W(y|x) + \log^2 \{P_n(x)W(y|x)\}). \end{aligned} \quad (440)$$

But then

$$0 \leq P_n(x)W(y|x) \left[ \log \frac{W(y|x)}{P_n W(y)} \right]^2 \quad (441)$$

$$\leq 2P_n(x)W(y|x)(\log^2 W(y|x) + \log^2 \{P_n(x)W(y|x)\}). \quad (442)$$

This is also true for  $P_n(x) = 0$  assuming the convention  $0 \log^2 0 = 0$ . Now continuity follows from the fact that  $x \log^2 \{\alpha x\}$  is continuous for  $x \in [0, 1]$  when defined as 0 for  $x = 0$ . Thus, continuity of  $U(P, W)$  is established.

To establish continuity of  $V(P, W)$  we are left to prove that

$$\sum_x P(x)D(W_x \| PW)^2$$

is continuous in  $P$ . Let us expand a single term here:

$$P(x) \left[ \sum_{y \in \mathcal{B}} W(y|x) \log \frac{W(y|x)}{PW(y)} \right]^2.$$

First notice that if  $P_n(x) \neq 0$  then continuity of this term follows from the fact that the argument of the logarithm is bounded away from 0 and  $+\infty$  for all  $y$  with  $W(y|x) > 0$ . So we are left with the case  $P_n(x) \rightarrow 0$ . To that end let us prove the inequality for  $P(x) > 0$ :

$$D(W_x \| PW) \leq 2H(W_x) + \log \frac{1}{P(x)}. \quad (443)$$

From here continuity follows as we can see that  $P_n(x)D(W_x \| P_n W)^2 \rightarrow 0$  because  $x \log x$  and  $x \log^2 x$  are continuous at zero.

We now prove inequality (443). From (439), we see that

$$\left| \log \frac{W(y|x)}{PW(y)} \right| \leq \log \frac{1}{W(y|x)} + \log \frac{1}{P(x)W(y|x)} \quad (444)$$

$$= 2 \log \frac{1}{W(y|x)} + \log \frac{1}{P(x)}. \quad (445)$$

Then,

$$D(W_x \| PW) \leq \sum_{x \in \mathcal{A}} W(y|x) \left| \log \frac{W(y|x)}{PW(y)} \right| \quad (446)$$

$$\leq 2H(W_x) + \log \frac{1}{P(x)}. \quad (447)$$

Thus,  $V(P, W)$  is continuous in  $P$ .

To establish continuity of  $T(P, W)$ , we again consider a single term:

$$P(x)W(y|x) \left| \log \frac{W(y|x)}{P_n W(y)} - D(W_x \| P_n W) \right|^3. \quad (448)$$

If  $W(y|x) = 0$  then this term is equal to zero regardless of  $P$ , and thus is continuous in  $P$ . Assume  $W(y|x) > 0$ . Take  $P_n \rightarrow P$ . If  $P(x) \neq 0$  then  $P_n W(y)$  is bounded away from 0 and thus  $\log \frac{W(y|x)}{P_n W(y)}$  tends to  $\log \frac{W(y|x)}{PW(y)}$ . Similarly, for any  $y'$  such that  $W(y'|x) > 0$  we have that  $P_n W(y')$  is also bounded away from 0. Thus,  $D(W_x \| P_n W)$  tends to  $D(W_x \| PW)$ .

We now assume that  $P_n(x) \rightarrow 0$  and must prove that (448) tends to 0. Using the inequality  $|a + b|^3 \leq 4(|a|^3 + |b|^3)$ , we obtain

$$P_n(x)W(y|x) \left| \log \frac{W(y|x)}{P_n W(y)} - D(W_x \| P_n W) \right|^3 \leq (449)$$

$$4P_n(x)W(y|x) \left| \log \frac{W(y|x)}{P_n W(y)} \right|^3 + 4P_n(x)W(y|x) D^3(W_x || P_n W). \quad (450)$$

Application of (443) immediately proves that the second term in the last inequality tends to zero. Continuity of the first term is established exactly like (438) with (440) replaced by

$$|\log W(y|x) - \log P_n W(y)|^3 \leq -4(\log^3 W(y|x) + \log^3(P_n(x)W(y|x))). \quad (451)$$

This proves continuity of  $T(P, W)$ .

Finally,  $V(P, W)$  and  $U(P, W)$  coincide on  $\Pi$  for the reason that, under any capacity-achieving distribution it is known that

$$D(W_x || PW) = \mathbb{E}[i(X; Y) | X = x] = C \quad P\text{-a.s.} \quad (452)$$

Indeed, then

$$U(P, W) \triangleq \mathbb{E}[(i - \mathbb{E} i)^2] \quad (453)$$

$$= \mathbb{E}[(i - C)^2] \quad (454)$$

$$= \mathbb{E}[(i - \mathbb{E}[i | X])^2] \quad (455)$$

$$= V(P, W). \quad (456)$$

To prove (436) consider the following chain of inequalities:

$$U(P, W) + [I(P, W)]^2 \triangleq \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P(x)W(y|x) \left[ \log^2 W(y|x) + \log^2 PW(y) - 2 \log W(y|x) \cdot \log PW(y) \right] \quad (457)$$

$$\leq \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P(x)W(y|x) [\log^2 W(y|x) + \log^2 PW(y)] \quad (459)$$

$$= \sum_{x \in \mathcal{A}} P(x) \left[ \sum_{y \in \mathcal{B}} W(y|x) \log^2 W(y|x) \right] \quad (460)$$

$$+ \left[ \sum_{y \in \mathcal{B}} PW(y) \log^2 PW(y) \right] \quad (461)$$

$$\leq \sum_{x \in \mathcal{A}} P(x)g(|\mathcal{B}|) + g(|\mathcal{B}|) \quad (462)$$

$$= 2g(|\mathcal{B}|) \quad (463)$$

where (459) is because  $2 \log W(y|x) \cdot \log PW(y)$  is always nonnegative, and (462) follows because each term in square-brackets can be upper-bounded using the following optimization problem:

$$g(n) \triangleq \sup_{a_j \geq 0: \sum_{j=1}^n a_j = 1} \sum_{j=1}^n a_j \log^2 a_j. \quad (464)$$

Since the  $x \log^2 x$  has unbounded derivative at the origin, the solution of (464) is always in the interior of  $[0, 1]^n$ . Then it is straightforward to show that for  $n > e$  the solution is actually  $a_j = \frac{1}{n}$ . For  $n = 2$  it can be found directly that  $g(2) = 0.5629 < 0.6$ . Finally, because of the symmetry, a similar argument can be made with  $|\mathcal{B}|$  replaced by  $|\mathcal{A}|$  and hence in (436) we are free to choose the best bound. ■

## APPENDIX F PROOF OF LEMMA 46

### A. Proof

Using Minkowski's inequality and the notation  $\|Z\|_p = (\mathbb{E} [|Z|^p])^{1/p}$ , we have

$$\|i(X; Y) - I(X; Y)\|_3 \leq \|i(X; Y)\|_3 + I(X; Y) \quad (465)$$

$$\leq \|\log W(Y|X) - \log P_Y(Y)\|_3 + I(X; Y) \quad (466)$$

$$\leq \|\log \frac{1}{W(Y|X)}\|_3 + \|\log \frac{1}{P_Y(Y)}\|_3 + I(X; Y) \quad (467)$$

$$\leq (|\mathcal{B}|27e^{-3} \log^3 e)^{1/3} + (|\mathcal{A}|27e^{-3} \log^3 e)^{1/3} + I(X; Y) \quad (468)$$

$$\leq (|\mathcal{A}|^{1/3} + |\mathcal{B}|^{1/3}) 3e^{-1} \log e + \log \min\{|\mathcal{A}|, |\mathcal{B}|\} \quad (469)$$

where (467) follows from  $x \log^3 \frac{1}{x} \leq (3e^{-1} \log e)^3$ . ■

## APPENDIX G

### A. Proof of Lemma 47

By Theorem 44 we have for any  $x$  and  $\delta$

$$\mathbb{P} \left[ x \leq \sum_{j=1}^n (Z_j - \mathbb{E} Z_j) < x + \delta \right] \quad (470)$$

$$\leq \int_{x/\sigma}^{(x+\delta)/\sigma} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt + \frac{12T}{\sigma^3} \quad (471)$$

$$\leq \left( \frac{\delta}{\sqrt{2\pi}} + \frac{12T}{\sigma^2} \right) \frac{1}{\sigma}. \quad (472)$$

On the other hand

$$\mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^n Z_j \right\} \mathbb{1}_{\left\{ \sum_{j=1}^n Z_j > A \right\}} \right] \quad (473)$$

$$\leq \sum_{l=0}^{\infty} \exp\{-A - l\delta\} \mathbb{P} \left[ A + l\delta \leq \sum_{j=1}^n Z_j < A + (l+1)\delta \right]. \quad (474)$$

Using (472) and  $\delta = \log 2$  we get (262) since

$$\sum_{l=0}^{\infty} 2^{-l} = 2. \quad (475)$$

#### APPENDIX H AN EXOTIC DMC

Let

$$W' = \begin{pmatrix} 1/3 & 0 & 0 & 1/3 & 2/13 \\ 0 & 1/3 & 0 & 1/3 & 2/13 \\ 0 & 0 & 1/3 & 1/3 & 2/13 \\ 1/3 & 1/3 & 1/3 & 0 & 3/13 \\ 1/3 & 1/3 & 1/3 & 0 & 4/13 \end{pmatrix}. \quad (476)$$

Now denote by  $x^*$  the unique negative root of the equation

$$(x-1) \left( \log \frac{13^7}{2833} - 7 \log(1-x) \right) + (6+7x) \log \frac{6+7x}{39} = -13 \log 3. \quad (477)$$

Then, replace the last column of  $W'$  with the column (only two decimals shown)

$$W'[0 \ 0 \ 0 \ x^* \ 1-x^*]^T = [0.04 \ 0.04 \ 0.04 \ 0.38 \ 0.50]^T. \quad (478)$$

The resulting channel matrix is

$$W = \begin{pmatrix} 1/3 & 0 & 0 & 1/3 & 0.04 \\ 0 & 1/3 & 0 & 1/3 & 0.04 \\ 0 & 0 & 1/3 & 1/3 & 0.04 \\ 1/3 & 1/3 & 1/3 & 0 & 0.38 \\ 1/3 & 1/3 & 1/3 & 0 & 0.50 \end{pmatrix}. \quad (479)$$

This matrix has full rank and so the capacity achieving distribution is unique. A simple observation shows that equiprobable  $P_Y^*$  is achievable by taking  $P_X^* = [1, 1, 1, 2, 0]/5$ . Finally, the conditional entropies  $H(Y|X=x)$  are all equal to  $\log 3$  as a consequence of the choice of  $x^*$ . It follows that  $P_X^*$  is the unique capacity achieving distribution. One can also check that  $V(P_X^*, W) = 0$  and  $V(W_5 || P_Y^*) > 0$ . So  $W$  is indeed an exotic channel. In fact, it can be shown that there is a sequence of distributions  $P_X^{(n)}$  such that Feinstein's lower bound for this channel exhibits  $nC + Fn^{1/3}$  behavior. Note that for an exotic channel and  $\epsilon > 1/2$  it is not optimal to choose  $P$  that achieves  $I(P, W) = C$  and  $U(P, W) = 0$  in Theorem 45, rather the optimal  $P$  will depend on  $n$ . The intuition behind this is that for small  $n$  it might be beneficial to choose  $P$  such that  $I(P, W) < C$  but  $U(P, W) > 0$  because for  $\epsilon > 1/2$  the  $\sqrt{n}$  term is positive and proportional to  $\sqrt{U(P, W)}$ .

This example has illustrated that the conditions for exotic channels are quite hard to satisfy (especially, making  $D(W_x || P_Y^*) = C$  but so that  $x$  does not participate in capacity achieving distributions); hence the name *exotic*.

#### APPENDIX I PROOF OF THEOREM 48

We must consider four cases separately:

1)  $\epsilon \leq 1/2$  and  $V_{\min} > 0$ ;

2)  $\epsilon \leq 1/2$  and  $V_{\min} = 0$ ;

3)  $\epsilon > 1/2$  and  $V_{\max} > 0$ ;

4)  $\epsilon > 1/2$  and  $V_{\max} = 0$ .

Compared to Strassen [31] we streamline the treatment of case 1 by using Lemma 64 and add the proofs for cases 3 and 4. The main idea for solving case 2 is due to Strassen. ■

The aim is to use Theorem 31 with  $F_n = T_{P_0}^n$ . To do so we need to select a distribution  $P_{Y^n}$  on  $\mathcal{A}^n$  and compute  $\inf_{x^n \in T_{P_0}^n} \beta_{1-\epsilon}^n(x^n, P_{Y^n})$ . Notice that the theorem is concerned only with codebooks over some fixed type. So, if  $P_{Y^n}$  is a product distribution then  $\beta_{1-\epsilon}^n(x^n, P_{Y^n})$  does not depend on  $x^n \in T_{P_0}^n$  and thus

$$\beta_{1-\epsilon}^n(x^n, P_{Y^n}) = \beta_{1-\epsilon}^n(P_{Y^n}). \quad (480)$$

For this reason, we will simply write  $\beta_{1-\epsilon}^n(P_{Y^n})$ , and even  $\beta_{1-\epsilon}^n$ , since the distribution  $P_{Y^n}$  will be apparent.

*Case 1:* Denote the closed  $\delta$ -neighborhood of the set of capacity-achieving distributions,  $\Pi$ , as

$$\Pi_\delta \triangleq \{P \in \mathcal{P} : d(P, \Pi) \leq \delta\}. \quad (481)$$

Here  $d(\cdot, \cdot)$  denotes Euclidean distance between vectors of  $\mathbb{R}^{|\mathcal{A}|}$ .

We fix some  $\delta > 0$  to be determined. First, we find  $\delta_1$  small enough so that everywhere on  $\Pi_{\delta_1}$  we have  $V(P, W) \geq V_{\min}/2$ . This is possible by the continuity of  $V(P, W)$ ; see Lemma 62 in Appendix E. Without loss of generality, we can assume that  $\mathcal{B}$  does not have inaccessible outputs, i.e., for every  $y_0 \in \mathcal{B}$  there is an  $x_0 \in \mathcal{A}$  such that  $W(y_0|x_0) > 0$ . Then, it is well known that for any  $P_1, P_2 \in \Pi$  the output distributions coincide, i.e.,  $P_1 W = P_2 W = P_Y^*$ , and also that this unique  $P_Y^*$  dominates all  $W(\cdot|x)$ . Since all outputs are accessible, this implies that  $P_Y^*(y) > 0$ ,  $y \in \mathcal{B}$ . Now for each  $y$ , the function  $PW(y)$  is linear in the input distribution  $P$ , and thus there is some  $\delta_2 > 0$  such that in the closed  $\delta_2$ -neighborhood of  $\Pi$  we have  $PW(y) > 0$  for all  $y \in \mathcal{B}$ . Set  $\delta = \min\{\delta_1, \delta_2\}$ . Fix  $n$  and  $P_0 \in \mathcal{P}_n$ . Choose the distribution  $P_{Y^n} = (P_0 W)^n$ , i.e.

$$P_{Y^n}(y^n) = \prod_{k=1}^n \sum_{a \in \mathcal{A}} P_0(a) W(y_k|a). \quad (482)$$

Then by Theorem 31 and the argument above, we have

$$\log M_{P_0}^*(n, \epsilon) \leq -\log \beta_{1-\epsilon}^n(x^n, P_{Y^n}) \quad (483)$$

where  $x^n$  is any element of  $T_{P_0}^n$ . The idea for lower bounding  $\beta_{1-\epsilon}^n$  is to apply Lemma 58 if  $P_0 \in \Pi_\delta$  and Lemma 59 (both in Appendix C) otherwise. In both cases,  $Q_i = P_{Y^i|X=x_i}$  and  $P_i = P_0 W$ . Note that there are  $nP_0(1)$  occurrences of  $P_{Y^i|X=1}$  among the  $Q_i$ 's,  $nP_0(2)$  occurrences of  $P_{Y^i|X=2}$ , etc. Thus, the quantities defined in Lemma 58 become

$$D_n = I(P_0, W) \quad (484)$$

$$V_n = V(P_0, W). \quad (485)$$

Suppose that  $P_0 \in \mathcal{P}_n \setminus \Pi_\delta$ ; then, applying Lemma 59 we obtain

$$\begin{aligned} \log M_{P_0}^*(n, \epsilon) &\leq -\log \beta_{1-\epsilon}^n \\ &\leq nI(P_0, W) + \sqrt{\frac{2nV(P_0, W)}{1-\epsilon}} + \log \frac{1-\epsilon}{2} \end{aligned} \quad (486)$$

$$\leq nC' + \sqrt{\frac{2M_V}{1-\epsilon}}\sqrt{n} + \log \frac{1-\epsilon}{2} \tag{487}$$

$$\tag{488}$$

where

$$C' = \sup_{P \in \mathcal{P} \setminus \Pi_\delta} I(P, W) < C \tag{489}$$

$$M_V = \max_{P \in \mathcal{P}} V(P, W) < \infty. \tag{490}$$

Since  $C' < C$  we can see that, even with  $F = 0$ , there exists  $N_1$  such that for all  $n \geq N_1$  the RHS of (488) is below the RHS of (272). So this proves (272) for  $P_0 \in \mathcal{P}_n \setminus \Pi_\delta$ . Now, consider  $P_0 \in \Pi_\delta$ . Recall that  $T_n$  in Lemma 58 is in fact

$$T_n = \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P_0(x)W(y|x) \left| \log \frac{W(y|x)}{P_0 W(y)} - D(W_x \| P_0 W) \right|^3 = T(P_0, W) \tag{491}$$

which, as shown in Lemma 62, is continuous on the compact set  $\mathcal{P}$  and thus has a finite upper bound:

$$T_n \leq M_T < \infty. \tag{492}$$

On the other hand, over  $\Pi_\delta$  we have  $V(P_0, W) \geq V_{\min}/2 > 0$ . In summary, we can upper bound  $B_n$  in Lemma 58 as

$$B_n = 6 \frac{T_n}{V_n^{3/2}} \leq M_B \triangleq \frac{6 \cdot 2^{3/2} M_T}{V_{\min}^{3/2}}. \tag{493}$$

Thus, we are ready to apply Lemma 58, namely to use (344) with  $\Delta = M_B - B_n + 1 \geq 1$  and to conclude that, for  $n$  sufficiently large

$$\log M_{P_0}^*(n, \epsilon) \leq nI(P_0, W) + \sqrt{nV(P_0, W)}Q^{-1} \left( 1 - \epsilon - \frac{M_B + 1}{\sqrt{n}} \right) + \frac{1}{2} \log n. \tag{494}$$

For  $n$  large, depending on  $M_B$ , we can expand  $Q^{-1}$  using Taylor's formula. In this way, we can conclude that there is a constant  $F_1$  such that

$$Q^{-1} \left( 1 - \epsilon - \frac{M_B + 1}{\sqrt{n}} \right) \leq Q^{-1}(1 - \epsilon) + \frac{F_1}{\sqrt{n}}. \tag{495}$$

Then for such  $n$  and a constant  $F_2$  (recall (490)) we have

$$\log M_{P_0}^*(n, \epsilon) \leq nI(P_0, W) + \sqrt{nV(P_0, W)}Q^{-1}(1 - \epsilon) + \frac{1}{2} \log n + F_2. \tag{496}$$

To conclude the proof we must maximize the RHS over  $P_0 \in \Pi_\delta$ . Note that this is the case treated in Lemmas 63 and 64. We want to use the latter one and need to check its conditions. From the definitions of  $I(P, W)$  and  $V(P, W)$  we can see that they are infinitely differentiable functions on  $\Pi_\delta$ . This is because all terms  $\log \frac{W(y|x)}{PW(y)}$  have arguments bounded away from

0 and  $+\infty$  by the choice of  $\Pi_\delta$ . Consequently, the conditions of Lemma 64 on  $g$  are automatically satisfied. We must now check the conditions on  $f$ . To that end, we can think of  $I(P, W)$  as a function of  $P$ , and write  $\nabla I(P)$  and  $\mathcal{H}(P)$  for the gradient vector and Hessian matrix correspondingly. To check the conditions on  $f$  in Lemma 64 it is sufficient to prove that for any  $P^* \in \Pi$ :

- 1)  $\ker \mathcal{H}(P^*) = \ker W$ , which is the set of all  $|\mathcal{A}|$ -vectors  $v$  such that  $\sum_{x \in \mathcal{A}} v(x)W(y|x) = 0$ ;
- 2) the largest nonzero eigenvalue of  $\mathcal{H}(P^*)$  is negative and bounded away from zero uniformly in the choice of  $P^* \in \Pi$ .

We first show why these two conditions are sufficient. It is known that  $\Pi$  consists of all distributions  $P$  that satisfy two conditions: 1)  $PW = P_Y^*$ ; and 2)  $P(x) > 0$  only when  $D(W_x \| P_Y^*) = C$ . Now take some  $P' \notin \Pi$  and denote by  $P^*$  the projection of  $P'$  onto a compact  $\Pi$ . Then write

$$P' = P^* + v = P^* + v_0 + v_\perp \tag{497}$$

where  $v_0$  is projection of  $v = (P' - P^*)$  onto  $\ker W$  and  $v_\perp$  is orthogonal to  $\ker W$ . Note that  $d(P', \Pi) = \|v\| \leq \delta$ . By Taylor's expansion we have

$$I(P') = I(P^*) + (v_0 + v_\perp)^T \nabla I(P^*) + \frac{1}{2} v_\perp^T \mathcal{H}(P^*) v_\perp + o(\|v\|^2). \tag{498}$$

Here we have used the fact that  $v_0^T \mathcal{H}(P^*) v_0 = 0$ . Since  $v_0 \in \ker W$  but  $P^* + \lambda v_0$  is not in  $\Pi$  for any  $\lambda > 0$ , we conclude that shifting along  $v_0$  must involve inputs with  $D(W_x \| P_Y^*) < C$ . But then  $I(P, W)$  decays linearly along this direction, i.e., there is some constant  $\Gamma_1 > 0$  such that

$$I(P^* + v_0) - I(P^*) = v_0^T \nabla I(P^*) \tag{499}$$

$$\leq -\Gamma_1 \|v_0\| \leq -\Gamma_1 \|v_0\|^2 \tag{500}$$

((500) assumes  $\delta \leq 1$ ). Then, substituting (500) into expansion for  $I(P')$  and upper bounding  $v_\perp^T \nabla I$  by zero we obtain

$$I(P') - I(P^*) \leq -\Gamma_1 \|v_0\|^2 - \frac{1}{2} \lambda \|v_\perp\|^2 + o(\|v\|^2) \tag{501}$$

where  $\lambda$  is the absolute value of the maximal nonzero eigenvalue of  $\mathcal{H}(P^*)$ . We will show that  $\lambda$  is uniformly bounded away from zero for any  $P^* \in \Pi$ . So we see that indeed  $I(P, W)$  decays not slower than quadratically in  $d(P, \Pi)$ .

Now we need to prove the assumed facts about the Hessian  $\mathcal{H}(P)$ . The differentiation can be performed without complications since on  $\Pi_\delta$  we always have  $PW(y) > 0$ . After some algebra we get

$$\mathcal{H}_{ij} \triangleq \frac{\partial^2 I(P)}{\partial P(i) \partial P(j)} = - \sum_{y \in \mathcal{B}} \frac{W(y|i)W(y|j)}{PW(y)}. \tag{502}$$

Thus, for any vector  $v$  we have

$$v^T \mathcal{H} v = \sum_{i,j} v_i \mathcal{H}_{ij} v_j \tag{503}$$

$$= - \sum_{y \in \mathcal{B}} \frac{\left( \sum_i v_i W(y|i) \right)^2}{PW(y)} \quad (504)$$

$$\leq - \frac{\|vW\|^2}{(PW)_{\max}} \quad (505)$$

where we have denoted formally  $vW = \sum_{x \in \mathcal{A}} v(x)W(y|x)$ , which is a vector of dimension  $|\mathcal{B}|$ . From (505) we can see that indeed  $v^T \mathcal{H} v = 0$  if and only if  $vW = 0$ . In addition, the maximal nonzero eigenvalue of  $\mathcal{H}(P)$  is always smaller than  $\frac{\lambda_{\min}(WW^T)}{(PW)_{\max}}$  for all  $P \in \Pi$ . Consequently, Lemma 64 applies to (496), and thus

$$\log M_{P_0}^*(n, \epsilon) \leq nC + \sqrt{nV_{\min}} Q^{-1}(1 - \epsilon) + \frac{1}{2} \log n + O(1). \quad (506)$$

This implies (272) if we note that  $Q^{-1}(1 - \epsilon) = -Q^{-1}(\epsilon)$ .

*Case 2:* The idea is to apply Theorem 31, but this time we fix the output distribution to be  $P_{Y^n} = (P_Y^*)^n$  for all types  $P_0$  (before we chose  $P_{Y^n} = (P_0 W)^n$  different for each type  $P_0$ ). It is well-known that

$$D(W \| P_Y^* | P_0) \leq D(W \| P_Y^* | P^*) = C. \quad (507)$$

This fact is crucial for proving the bound.

Note that  $V(W_x \| P_Y^*)$  is defined and finite since all  $W_x \ll P_Y^*$ . Denote a special subset of *nonzero-variance inputs* as

$$\mathcal{A}_+ \triangleq \{x \in \mathcal{A} : V(W_x \| P_Y^*) > 0\}. \quad (508)$$

And also for every  $P_0 \in \mathcal{P}_n$  denote  $m(P_0) = nP_0(\mathcal{A}_+)$  which is the number of nonzero-variance letters in any  $x \in T_{P_0}^n$ . Also note that there are minimal and maximal variances  $V_M \geq V_m > 0$  such that  $V_m \leq V(W_x \| P_Y^*) \leq V_M$  for all  $x \in \mathcal{A}_+$ .

Since  $P_{Y^n}$  is a product distribution

$$\log M_{P_0}^*(n, \epsilon) \leq -\log \beta_{1-\epsilon}^n(x^n, P_{Y^n}) \quad (509)$$

for all  $x^n \in T_{P_0}^n$ . We are going to apply Lemmas 58 and 59, Appendix C, and so need to compute  $D_n$ ,  $V_n$  and an upper bound on  $B_n$ . We have

$$D_n = D(W \| P_Y^* | P_0) \quad (510)$$

$$V_n = V(W \| P_Y^* | P_0). \quad (511)$$

To upper bound  $B_n$  we must lower bound  $V_n$  and upper bound  $T_n$ . Note that

$$V(W \| P_Y^* | P_0) \geq \frac{m(P_0)}{n} V_m. \quad (512)$$

For  $T_n$ , we can write

$$\begin{aligned} T_n &\triangleq \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} P_0(x) W(y|x) \left| \log \frac{W(y|x)}{P_Y^*(y)} - D(W_x \| P_Y^*) \right|^3 \\ &= \sum_{x \in \mathcal{A}} P_0(x) T(x). \end{aligned} \quad (513)$$

Here, the  $T(x)$ 's are all finite and  $T(x) = 0$  iff  $x \notin \mathcal{A}_+$ . Thus, for  $x \in \mathcal{A}_+$  there is one maximal  $T^* = \max_{x \in \mathcal{A}} T(x)$ , and we have

$$T_n \leq \frac{m(P_0)}{n} T^*. \quad (514)$$

Then, we see that

$$B_n \triangleq \frac{T_n}{V_n^{3/2}} \leq \sqrt{\frac{n}{m(P_0)}} \frac{T^*}{V_m^{3/2}} \triangleq \sqrt{\frac{n}{m(P_0)}} M_B. \quad (515)$$

So we apply Lemma 58 with

$$\Delta = \sqrt{\frac{n}{m(P_0)}} (M_B + 1) - B_n \geq \sqrt{\frac{n}{m(P_0)}} \geq 1. \quad (516)$$

Using (344) and lower bounding  $\log \Delta$  via (516) we have

$$\begin{aligned} \log \beta_{1-\epsilon}^n &\geq -nD(W \| P_Y^* | P_0) \\ &\quad - \sqrt{nV(W \| P_Y^* | P_0)} Q^{-1} \left( 1 - \epsilon - \frac{M_B + 1}{\sqrt{m(P_0)}} \right) \\ &\quad - \frac{1}{2} \log n. \end{aligned} \quad (517)$$

Now, it is an elementary analytical fact that it is possible to choose a  $\delta_0 < 1 - \epsilon$  and  $\Gamma_2 > 0$  such that

$$Q^{-1}(1 - \epsilon - z) \leq Q^{-1}(1 - \epsilon) + \Gamma_2 z \forall z \in [0, \delta_0]. \quad (518)$$

We now split types in  $\mathcal{P}_n$  into two classes,  $\mathcal{P}_A$  and  $\mathcal{P}_B$

$$P_0 \in \mathcal{P}_A \iff m(P_0) \geq m_* \quad \mathcal{P}_B = \mathcal{P}_n \setminus \mathcal{P}_A. \quad (519)$$

Here  $m_*$  is chosen so that  $\frac{M_B + 1}{\sqrt{m_*}} \leq \delta_0$ . Then, for all types in  $\mathcal{P}_A$ , we have

$$Q^{-1} \left( 1 - \epsilon - \frac{M_B + 1}{\sqrt{m(P_0)}} \right) \leq Q^{-1}(1 - \epsilon) + \frac{\Gamma_3}{\sqrt{m(P_0)}}. \quad (520)$$

Notice also that with this choice of  $x_0$  and  $m_*$ , the argument of  $Q^{-1}$  in (517) is positive and the bound is applicable to all types in  $\mathcal{P}_A$ . Substituting (507) we have, for any  $P_0 \in \mathcal{P}_A$

$$\begin{aligned} \log \beta_{1-\epsilon}^n &\geq -nC - \sqrt{nV(W \| P_Y^* | P_0)} Q^{-1}(1 - \epsilon) \\ &\quad - \Gamma_3 \sqrt{\frac{nV(W \| P_Y^* | P_0)}{m(P_0)}} - \frac{1}{2} \log n. \end{aligned} \quad (521)$$

Now notice that  $Q^{-1}(1 - \epsilon) \leq 0$  (this is the key difference with Case 4) and also that

$$V(W \| P_Y^* | P_0) \leq \frac{m(P_0)}{n} V_M. \quad (522)$$

Finally, for  $P_0 \in \mathcal{P}_A$  we have

$$\log M_{P_0}^*(n, \epsilon) \leq nC + \Gamma_3 \sqrt{V_M} + \frac{1}{2} \log n. \quad (523)$$

Now for types in  $\mathcal{P}_B$  we have  $m(P_0) < m_*$  and thus

$$nV(W \| P_Y^* | P_0) \leq m_* V_M. \quad (524)$$

So Lemma 59 yields

$$\log M_{P_0}^*(n, \epsilon) \leq nC + \sqrt{\frac{2m_* V_M}{1 - \epsilon}} - \log \frac{1 - \epsilon}{2}. \quad (525)$$

In summary, we see that in both cases,  $\mathcal{P}_A$  and  $\mathcal{P}_B$ , inequalities (523) and (525) imply (272) for  $n \geq 1$ .

*Case 3:* The proof for this case is analogous to that for Case 1, except that when applying Lemma 64 we must choose  $g^* = \sqrt{V_{\max}}$  because the sign of  $Q^{-1}(1 - \epsilon)$  is positive this time. An additional difficulty is that it might be possible that  $V_{\max} > 0$  but  $V_{\min} = 0$ . In this case the bound (493) is no longer applicable. What needs to be done is to eliminate types inside  $\Pi_\delta$  with small variance

$$\Pi_V = \{P \in \Pi_\delta : V(P, W) < A\} \quad (526)$$

where

$$A < \frac{1 - \epsilon}{2} V_{\max} (Q^{-1}(\epsilon))^2. \quad (527)$$

Then, for types in  $\Pi_V$  we can apply the fixed-blocklength bound in Lemma 59. For the remaining types in  $\Pi_\delta \setminus \Pi_V$  the argument in Case 1 works, after  $V_{\min}$  is replaced by  $A$  in (493).

*Case 4:* Fix a type  $P_0 \in \mathcal{P}_n$  and use  $P_{Y^n} = \prod_1^n(P_0 W)$ . Then, a similar argument to that for Case 2 and Lemma 59 yields

$$\log M_{P_0}^*(n, \epsilon) \leq nI(P_0, W) + \sqrt{\frac{2nV(P_0, W)}{1 - \epsilon}} + \log \frac{1 - \epsilon}{2} \quad (528)$$

for all  $n \geq 1$ . We need to maximize the RHS of this bound over  $P_0 \in \mathcal{P}$ . This can be done similarly to Lemma 64. The problem here, however, is that  $V(P, W) = 0$  for  $P \in \Pi$ . Thus, even though  $V(P, W)$  is differentiable in some neighborhood of  $\Pi$ ,  $\sqrt{V(P, W)}$  is not. This is how a term of order  $n^{1/3}$  can appear. Indeed, suppose that there is some direction  $v$  along which  $I(P + \lambda v)$  decays quadratically, while  $V(P + \lambda v)$  is linear. For example

$$I(P + \lambda v) = C - \Gamma_4 \lambda^2 + o(\lambda^2) \quad (529)$$

$$V(P + \lambda v) = \Gamma_5 \lambda + o(\lambda). \quad (530)$$

Then it is not hard to see that

$$\begin{aligned} \max_\lambda \left\{ nI(P + \lambda v) + \sqrt{nV(P + \lambda v)} \right\} \\ = nC + \Gamma_6 n^{1/3} + o(n^{1/3}). \end{aligned} \quad (531)$$

Such a direction can only exist if all the conditions of the exotic DMC are satisfied. This can be proved by computing gradients of  $I(P, W)$  and  $V(P, W)$ .

#### APPENDIX J

##### AUXILIARY MAXIMIZATION LEMMAS

This appendix is concerned with the behavior of the maximum of  $nf(x) + \sqrt{ng(x)}$  for large  $n$ , for arbitrary continuous  $f$  and  $g$ .

*Lemma 63:* Let  $D$  be a compact metric space. Suppose  $f : D \mapsto \mathbb{R}$  and  $g : D \mapsto \mathbb{R}$  are continuous. Define

$$f^* = \max_{x \in D} f(x) \quad (532)$$

and

$$g^* = \sup_{\{x: f(x)=f^*\}} g(x). \quad (533)$$

Then,

$$\max_{x \in D} [nf(x) + \sqrt{ng(x)}] = nf^* + \sqrt{ng^*} + o(\sqrt{n}). \quad (534)$$

The message of this lemma is that, for continuous  $f$  and  $g$ ,

$$\max_x \{nf(x) + \sqrt{ng(x)}\} \approx nf(x^{**}) + \sqrt{ng(x^{**})} \quad (535)$$

where  $x^{**}$  is found by first maximizing  $f(x)$  and then maximizing  $g(x)$  over the set of maximizers of  $f(x)$ .

If we assume more about  $f$  and  $g$ , then a stronger result can be stated. The assumptions below essentially mean that  $f$  is twice differentiable near  $f^*$  with negative-definite Hessian and  $g$  is differentiable. Without such assumptions Lemma 63 appears to be the best possible result; see the example after the proof of Lemma 63 below.

*Lemma 64:* In the notation of previous lemma, denote

$$D_0 \triangleq \{x : f(x) = f^*\} \quad (536)$$

$$D_\delta \triangleq \{x : d(x, D_0) \leq \delta\} \quad (537)$$

where  $d(\cdot, \cdot)$  is a metric. Suppose that for some  $\delta > 0$  and some constants  $f_1 > 0$  and  $f_2$  we have

$$f(x) - f^* \leq -f_1 d(x, D_0)^2 \quad (538)$$

$$|g(x) - g^*| \leq f_2 d(x, D_0) \quad (539)$$

for all  $x \in D_\delta$ . Then

$$\max_{x \in D} [nf(x) + \sqrt{ng(x)}] = nf^* + \sqrt{ng^*} + O(1). \quad (540)$$

*Proof of Lemma 63:* Denote

$$F(x, n) = nf(x) + \sqrt{ng(x)} \quad (541)$$

$$F^*(n) = \max_{x \in D} F(x, n). \quad (542)$$

Then (534) is equivalent to a pair of statements:

$$\lim_{n \rightarrow \infty} \frac{1}{n} F^*(n) = f^* \quad (543)$$

$$\lim_{n \rightarrow \infty} \frac{F^*(n) - nf^*}{\sqrt{n}} = g^* \quad (544)$$

which we are going to prove. First we note that because of the compactness of  $D$  both  $f$  and  $g$  are bounded. Now

$$F(x, n) \leq nf^* + \sqrt{ng_{\max}} \quad (545)$$

which implies

$$\frac{1}{n} F^*(n) \leq f^* + \frac{1}{\sqrt{n}} g_{\max} \quad (546)$$

which in turn implies

$$\limsup \frac{1}{n} F^*(n) \leq f^*. \quad (547)$$

On the other hand, if we take  $x^*$  to be any  $x \in D$  maximizing  $f(x)$  then

$$F^*(n) = \max_x F(x, n) \geq F(x^*, n) = n f^* + \sqrt{n} g(x^*). \quad (548)$$

Thus

$$\liminf \frac{1}{n} F^*(n) \geq f^* \quad (549)$$

and the first statement is proved. Now define

$$D_1 = \{x \in D : f(x) = f^*\} \quad (550)$$

which is also compact. Thus, there exists an (possibly nonunique)  $x^{**}$  maximizing  $g(x)$  on  $D_0$ :

$$x^{**} = \operatorname{argmax}_{x \in D_0} g(x) \quad (551)$$

$$g(x^{**}) = g^*. \quad (552)$$

By definition

$$F^*(n) - n f^* \geq F(x^{**}, n) - n f^* = \sqrt{n} g^*. \quad (553)$$

Thus

$$\liminf \frac{F^*(n) - n f^*}{\sqrt{n}} \geq g^*. \quad (554)$$

On the other hand,  $F(x, n)$  is continuous on  $D$ , so that

$$F^*(n) = F(x_n^*, n). \quad (555)$$

Then notice that

$$F^*(n) - n f^* = n(f(x_n^*) - f^*) + \sqrt{n} g(x_n^*) \quad (556)$$

$$\leq \sqrt{n} g(x_n^*) \quad (557)$$

where the last inequality follows because  $f(x_n^*) \leq f^*$ . Now we see that

$$\frac{F^*(n) - n f^*}{\sqrt{n}} \leq g(x_n^*). \quad (558)$$

Denoting

$$\phi(n) \triangleq \frac{F^*(n) - n f^*}{\sqrt{n}} \quad (559)$$

there exists a sequence  $\{n_k\}$  such that

$$\phi(n_k) \rightarrow \limsup \phi(n) \quad \text{as } k \rightarrow \infty. \quad (560)$$

For that sequence we have

$$\phi(n_k) \leq g(x_{n_k}^*). \quad (561)$$

Since the  $x_{n_k}^*$ 's all lie in the compact  $D$ , there exists a convergent subsequence<sup>20</sup>:

$$y_l \triangleq x_{n_{k_l}}^* \rightarrow x_0. \quad (562)$$

We will now argue that  $f(x_0) = f^*$ . As we have just shown,

$$\frac{1}{n_{k_l}} F^*(n_{k_l}) \rightarrow f^* \quad (563)$$

where

$$F^*(n_{k_l}) = F(y_l, n_{k_l}) = n_{k_l} f(y_l) + \sqrt{n_{k_l}} g(y_l). \quad (564)$$

Thus, since  $g(x)$  is bounded

$$\lim_{l \rightarrow \infty} \frac{1}{n_{k_l}} F^*(n_{k_l}) = \lim_{l \rightarrow \infty} f(y_l) = f(x_0) \quad (565)$$

where the last step follows from the continuity of  $f$ . So indeed

$$f(x_0) = f^* \iff x_0 \in D_0 \implies g(x_0) \leq g^*. \quad (566)$$

Now we recall that

$$\phi(n_{k_l}) \leq g(y_l) \quad (567)$$

and by taking the limit as  $l \rightarrow \infty$ , we obtain

$$\limsup \phi(n) = \lim_{l \rightarrow \infty} \phi(n_{k_l}) \quad (568)$$

$$\leq \lim_{l \rightarrow \infty} g(y_l) = g(x_0) \quad (569)$$

$$\leq g^*. \quad (570)$$

So we have shown

$$\lim \frac{F^*(n) - n f^*}{\sqrt{n}} = g^*. \quad (571)$$

■

Lemma 63 is tight in the sense that term  $o(\sqrt{n})$  cannot be improved without further assumptions. Indeed, take  $f(x) = -x^2$  and  $g(x) = x^{1/k}$  for some  $k \in \mathbb{Z}_+$  on  $[-1, 1]$ . Then, a simple calculation shows that

$$\max_{x \in [-1, 1]} \{n f(x) + \sqrt{n} g(x)\} = \text{const} \cdot n^{\frac{k-1}{2k-1}} \quad (572)$$

and the power of  $n$  can be arbitrary close to  $\sqrt{n}$ .

Lemma 63 can be generalized to any finite set of ‘‘basis terms’’, instead of  $\{n, \sqrt{n}\}$ . In this case, the only requirement would be that  $u_{j+1}(n) = o(u_j(n))$ .

<sup>20</sup>This is the only place where we use the metric-space nature of  $D$ . Namely we need sequential compactness to follow from compactness. Thus, Lemma 63 holds in complete generality for an arbitrary compact topological space  $D$  that is first-countable (i.e., every point has a countable neighborhood basis).



*Proof of Lemma 64:* Because of the boundedness of  $g(x)$ , the points  $x_n^*$  must all lie in  $D_\delta$  for  $n$  sufficiently large. So, for such  $n$  we have

$$\max_{x \in D} F(x, n) = \max_{x \in D_\delta} F(x, n) \quad (573)$$

$$\begin{aligned} \max_{x \in D_\delta} F(x, n) &= n f^* + \sqrt{n} g^* \\ &+ [n(f(x_n^*) - f^*) + \sqrt{n}(g(x_n^*) - g^*)]. \end{aligned} \quad (574)$$

We can now bound the term in brackets by using conditions in the lemma:

$$0 \leq [n(f(x_n^*) - f^*) + \sqrt{n}(g(x_n^*) - g^*)] \quad (575)$$

$$\leq -f_1 (\sqrt{n}d(x_n^*, D_0))^2 + f_2 (\sqrt{n}d(x_n^*, D_0)). \quad (576)$$

Now we see that we have a quadratic polynomial in the variable  $y \triangleq \sqrt{n}d(x_n^*, D_0)$ . Since  $f_1 > 0$  it has a maximum equal to  $\frac{f_2^2}{4f_1^2}$ . Then

$$0 \leq [n(f(x_n^*) - f^*) + \sqrt{n}(g(x_n^*) - g^*)] \leq \frac{f_2^2}{4f_1^2} \quad (577)$$

and we see that residual term is  $O(1)$ . This establishes (540). ■

## APPENDIX K

### REFINED EXPANSIONS FOR THE BSC AND BEC

#### A. Proof of Theorem 52

The converse bound was computed in Section III-H in (180) and (173). To analyze the asymptotics of  $\beta_\alpha^n$  we proceed as in the proof of Theorem 48, Case 1. Similarly to (496), we obtain

$$\begin{aligned} \log_2 M^*(n, \epsilon) &\leq n(1 - h(\delta)) \\ &- \sqrt{n\delta(1 - \delta)} \log_2 \frac{1 - \delta}{\delta} Q^{-1}(\epsilon) + \frac{1}{2} \log_2 n + O(1). \end{aligned} \quad (578)$$

Note that because of Theorem 28 this upper bound holds even if  $\epsilon$  is an average probability of error.

We now return to the achievability part. In order to obtain the constant in the  $\log n$  term we use Theorem 33, as none of the other bounds is tight enough to yield the right  $\log n$  term. First, denote

$$S_n^k \triangleq 2^{-n} \sum_{l=0}^k \binom{n}{l}. \quad (579)$$

Then (162) implies the existence of an  $(n, M, \epsilon)$  code (maximal probability of error) with

$$\epsilon \leq \sum_{k=0}^n \binom{n}{k} \delta^k (1 - \delta)^{n-k} \min \{1, MS_n^k\}. \quad (580)$$

We will argue that (580) implies a lower bound on  $M^*$  with a matching  $\log n$  term.

Without loss of generality, assume  $\delta < 1/2$ ; choose any  $r \in (\delta, 1/2)$  and set

$$q = \frac{r}{1 - r} < 1 \quad (581)$$

$$K = n\delta + \sqrt{n\delta(1 - \delta)}Q^{-1} \left( \epsilon - \frac{B + G}{\sqrt{n}} \right) \quad (582)$$

where  $B$  denotes the Berry–Esseen constant for a binomial  $(n, \delta)$  distribution

$$G = \frac{G_1}{(1 - q)^2} \quad (583)$$

and  $G_1$  is a constant (guaranteed to exist according [46]) such that for all  $k = 0, \dots, n$

$$\binom{n}{k} \delta^k (1 - \delta)^{n-k} \leq \frac{G_1}{\sqrt{n}}. \quad (584)$$

Then from Berry–Esseen Theorem 44 we obtain

$$\sum_{k > K} \binom{n}{k} \delta^k (1 - \delta)^{n-k} \leq \epsilon - \frac{G}{\sqrt{n}}. \quad (585)$$

It is also clear that for all sufficiently large  $n$  we have  $K < rn$ . Now, observe the following inequality, valid for  $k = 1, \dots, n - 1$  and  $j = -(n - k), \dots, k$ :

$$\binom{n}{k - j} \leq \binom{n}{k} \left( \frac{k}{n - k} \right)^j. \quad (586)$$

Consider any  $M$  such that  $MS_n^K \leq 1$ , then

$$M \sum_{k=0}^K S_n^k = M \sum_{t=0}^K (K - t + 1) \binom{n}{t} 2^{-n} \quad (587)$$

$$= M \sum_{l=0}^K (l + 1) \binom{n}{K - l} 2^{-n} \quad (588)$$

$$\leq MS_n^K \sum_{l=0}^K (l + 1) \left( \frac{K}{n - K} \right)^l \quad (589)$$

$$\leq MS_n^K \sum_{l=0}^K (l + 1) q^l \quad (590)$$

$$\leq MS_n^K \sum_{l=0}^{\infty} (l + 1) q^l \quad (591)$$

$$\leq \frac{1}{(1 - q)^2}. \quad (592)$$

If  $MS_n^K \leq 1$  then by (592)

$$\sum_{k=0}^K \binom{n}{k} \delta^k (1 - \delta)^{n-k} MS_n^k \leq \frac{G}{\sqrt{n}}. \quad (593)$$

We can now see that (580) implies that

$$M^*(n, \epsilon) \geq \frac{1}{S_n^K}. \quad (594)$$

Indeed, choose  $M = \frac{1}{S_n^k}$ . Then from (585) and (593) it follows that

$$\sum_{k=0}^n \binom{n}{k} \delta^k (1-\delta)^{n-k} \min\{1, MS_n^k\} \quad (595)$$

$$\leq \sum_{k=0}^K \binom{n}{k} \delta^k (1-\delta)^{n-k} + \sum_{k>K} \binom{n}{k} \delta^k (1-\delta)^{n-k} \quad (596)$$

$$\leq \frac{G}{\sqrt{n}} + \epsilon - \frac{G}{\sqrt{n}} \quad (597)$$

$$= \epsilon. \quad (598)$$

Finally, we must upper bound  $\log S_n^K$  up to  $O(1)$  terms. This is simply an application of (586):

$$S_n^K = 2^{-n} \sum_{k=0}^K \binom{n}{k} \quad (599)$$

$$\leq 2^{-n} \binom{n}{K} \sum_{l=0}^{\infty} \left(\frac{K}{n-K}\right)^l \quad (600)$$

$$\leq 2^{-n} \binom{n}{K} \frac{n-K}{n-2K}. \quad (601)$$

For  $n$  sufficiently large  $n-2K$  will become larger than  $n(1-2r)$ ; thus for such  $n$  we have  $\frac{n-K}{n-2K} \leq \frac{1}{1-2r}$  and hence

$$\log_2 S_n^K \leq -n + \log_2 \binom{n}{K} + O(1). \quad (602)$$

Using Stirling's approximation we obtain the inequality

$$\binom{n}{K} \leq \frac{e^{1/12}}{\sqrt{2\pi}} \sqrt{\frac{n}{K(n-K)}} \exp(nh(K/n)). \quad (603)$$

Substituting  $K$  from (582) and applying Taylor's formula to  $h(p)$  implies

$$\begin{aligned} \log_2 S_n^K &\leq n(h(\delta) - 1) \\ &\quad + \sqrt{n\delta(1-\delta)} \log_2 \frac{1-\delta}{\delta} Q^{-1} \left( \epsilon - \frac{B+G}{\sqrt{n}} \right) \\ &\quad - \frac{1}{2} \log_2 n + O(1). \end{aligned} \quad (604)$$

Finally, applying Taylor's formula to  $Q^{-1}$ , we conclude

$$\begin{aligned} \log_2 S_n^K &\leq n(h(\delta) - 1) \\ &\quad + \sqrt{n\delta(1-\delta)} \log_2 \frac{1-\delta}{\delta} Q^{-1}(\epsilon) - \frac{1}{2} \log_2 n + O(1). \end{aligned} \quad (605)$$

Substituting this into (594) we obtain the sought-after expansion.  $\blacksquare$

### B. Proof of Theorem 53

The achievability part of (290) is established by (276). The converse in Theorem 48 yields the wrong  $\log n$  term; instead, we use the stronger converse in Theorem 38 (which holds for average error probability). Since any  $(n, M, \epsilon)$  code must satisfy this bound then we must simply find  $M$  so large that the left-hand side (LHS) is larger than a given  $\epsilon$ . We can then conclude

that  $M^*(n, \epsilon)$  is upper bounded by such  $M$ . We observe that by (584)

$$\sum_{\ell=\lfloor n-\log_2 M \rfloor + 1}^n \binom{n}{\ell} \delta^\ell (1-\delta)^{n-\ell} 2^{n-\ell-\log_2 M} \leq \frac{2G_1}{\sqrt{n}}. \quad (606)$$

Then, denote by  $B$  the usual Berry–Esseen constant for a binomial distribution, and set

$$\log_2 M = n(1-\delta) - \sqrt{n\delta(1-\delta)} Q^{-1} \left( \epsilon + \frac{B+2G_1}{\sqrt{n}} \right). \quad (607)$$

Then from Berry–Esseen Theorem 44, we obtain

$$\sum_{l \geq n - \log_2 M} \binom{n}{l} \delta^l (1-\delta)^{n-l} \geq \epsilon + \frac{2G_1}{\sqrt{n}}. \quad (608)$$

Finally from (606) we conclude that

$$\sum_{\ell=\lfloor n-\log_2 M \rfloor + 1}^n \binom{n}{\ell} \delta^\ell (1-\delta)^{n-\ell} (1 - 2^{n-\ell-\log_2 M}) \geq \epsilon \quad (609)$$

and hence

$$\begin{aligned} \log_2 M^*(n, \epsilon) &\leq n(1-\delta) - \sqrt{n\delta(1-\delta)} Q^{-1} \left( \epsilon + \frac{B+2G_1}{\sqrt{n}} \right) \quad (610) \\ &= n(1-\delta) - \sqrt{n\delta(1-\delta)} Q^{-1}(\epsilon) + O(1) \quad (611) \end{aligned}$$

where (611) follows from Taylor's formula.  $\blacksquare$

## APPENDIX L PROOF OF THEOREM 54

It is convenient to split the proof of Theorem 54 into three parts. We first address the converse parts.

*Theorem 65:* For the AWGN channel with SNR  $P$  and  $\epsilon \in (0, 1)$  and equal-power constraint we have

$$\log M_e^*(n, \epsilon, P) \leq nC - \sqrt{nV} Q^{-1}(\epsilon) + \frac{1}{2} \log n + O(1) \quad (612)$$

where the capacity  $C$  and dispersion  $V$  are defined in (292) and (293).

*Proof:* Take  $A, B, P_{Y^n}$  and  $P_{Y^n|X^n=x^n}$  as in Section III-J. There we have shown that, for any  $x^n \in \mathbb{F}_n$  the distribution of  $i(x^n; Y^n)$  is the same as that of  $H_n$  in (205). Thus, using (106), we have for any  $\zeta_n$ ,

$$\inf_{x^n \in \mathbb{F}_n} \beta_{1-\epsilon}^n(x^n) = \beta_{1-\epsilon}^n(x^n) \quad (613)$$

$$\geq \exp(\zeta_n - nC) \left( 1 - \epsilon - \mathbb{P} \left[ \sum_{i=1}^n S_i \leq \zeta_n \right] \right). \quad (614)$$

with

$$S_i = \frac{P \log e}{2(1+P)} \left( Z_i^2 - 2 \frac{Z_i}{\sqrt{P}} - 1 \right) \quad (615)$$

and the  $Z_i$ 's are i.i.d. standard normal. Note that  $\mathbb{E}[S_i] = 0$  and

$$\text{Var}(S_i) = \left( \frac{P \log e}{2(1+P)} \right)^2 \mathbb{E} \left[ \left( Z_i^2 - 2 \frac{Z_i}{\sqrt{P}} - 1 \right)^2 \right] \quad (616)$$

$$= \left( \frac{P \log e}{2(1+P)} \right)^2 \mathbb{E} \left[ Z_i^4 + \left( \frac{4}{P} - 2 \right) Z_i^2 + 1 \right] \quad (617)$$

$$= \left( \frac{P \log e}{2(1+P)} \right)^2 \left( \frac{4}{P} + 2 \right) \quad (618)$$

$$= V. \quad (619)$$

Furthermore, define

$$B(P) = \frac{6\mathbb{E}[|S_i|^3]}{V^{3/2}} \quad (620)$$

$$N_c(P, \epsilon) = \left( \frac{2B(P)}{1-\epsilon} \right)^2. \quad (621)$$

Then for  $n > N_c(P, \epsilon)$  we have

$$\alpha_n = 1 - \epsilon - \frac{2B(P)}{\sqrt{n}} > 0. \quad (622)$$

For such  $n$  choose

$$\zeta_n = -\sqrt{nV}Q^{-1}(\alpha_n). \quad (623)$$

Then from Theorem 44, we have

$$\mathbb{P} \left[ \sum_{i=1}^n S_i \leq \zeta_n \right] \leq \alpha_n + \frac{B(P)}{\sqrt{n}} \quad (624)$$

$$\leq 1 - \epsilon - \frac{B(P)}{\sqrt{n}}. \quad (625)$$

On substituting (625) into (613) we obtain

$$\beta_{1-\epsilon}^n \geq \exp(\zeta_n - nC) \frac{B(P)}{\sqrt{n}}. \quad (626)$$

Using Theorem 31, this implies

$$\log M_e^*(n, \epsilon, P) \leq nC - \zeta_n + \frac{1}{2} \log n - \log B(P). \quad (627)$$

From Taylor's theorem, for some  $\theta \in \left[ 1 - \epsilon - \frac{2B(P)}{\sqrt{n}}, 1 - \epsilon \right]$ , we have

$$\zeta_n = -\sqrt{nV}Q^{-1}(1 - \epsilon) + 2B(P)\sqrt{V} \frac{dQ^{-1}}{dx}(\theta). \quad (628)$$

Without loss of generality, we assume that  $0 < \epsilon < 1 - \frac{2B(P)}{\sqrt{n}}$ , for all  $n > N_c(P, \epsilon)$  (otherwise just increase  $N_c(P, \epsilon)$  until this is true). Since  $\frac{d}{dx}Q^{-1}$  is a continuous function on  $(0, 1)$ , we can lower bound  $\frac{d}{dx}Q^{-1}(\theta)$  by

$$g_1(P, \epsilon) = \min_{\theta \in [\alpha_1, 1-\epsilon]} \frac{d}{dx}Q^{-1}(\theta) \quad (629)$$

where  $\alpha_1 = 1 - \epsilon - \frac{2B(P)}{\sqrt{N_c(P, \epsilon+1)}}$ . Note that  $g_1(P, \epsilon)$  is a continuous function of  $P$  and  $\epsilon$ . This results in

$$\zeta_n \geq -\sqrt{nV}Q^{-1}(1 - \epsilon) + g_1(P, \epsilon)2B(P)\sqrt{V}. \quad (630)$$

Substituting this bound into (627) and defining

$$g_c(P, \epsilon) = -2B(P)\sqrt{V}g_1(P, \epsilon) - \log B(P) \quad (631)$$

we arrive at

$$\log M_e^*(n, \epsilon, P) \leq nC + \sqrt{nV}Q^{-1}(1 - \epsilon) + \frac{1}{2} \log n + g_c(P, \epsilon). \quad (632)$$

*Corollary 66:* For the AWGN channel with SNR  $P$  and for each  $0 < \epsilon < 1$ , we have

$$M_m^*(n, \epsilon, P) \leq nC - \sqrt{nV}Q^{-1}(\epsilon) + \frac{1}{2} \log n + O(1) \quad (633)$$

$$M_a^*(n, \epsilon, P) \leq nC - \sqrt{nV}Q^{-1}(\epsilon) + \frac{3}{2} \log n + O(1). \quad (634)$$

*Proof:* By Lemma 39 we have

$$\log M_m^*(n, \epsilon, P) \leq \log M_e^*(n + 1, \epsilon, P). \quad (635)$$

Therefore from (612) and Taylor's theorem we get (633).

To prove (634) we set

$$N(\epsilon, P) = \max_{P_1 \in [P, 2P]} N_c(\epsilon, P_1) \quad (636)$$

$$g(\epsilon, P) = \max_{P_1 \in [P, 2P]} g_c(\epsilon, P_1) \quad (637)$$

where  $N_c$  and  $g_c$  are continuous functions defined in (621) and (631). Now set  $P_n = (1 + 1/n)P$  and use Lemma 39. Then for all  $n > N(\epsilon, P)$  according to Theorem 65 we have

$$\log M_a^*(n, \epsilon, P) \leq -\log \left( 1 - \frac{P}{P_n} \right) + \log M_m^*(n, \epsilon, P_n) \quad (638)$$

$$\leq \log(n + 1) + \log M_e^*(n + 1, \epsilon, P_n) \quad (639)$$

$$\leq (n + 1)C(P_n) - \sqrt{(n + 1)V(P_n)}Q^{-1}(\epsilon) + \frac{3}{2} \log(n + 1) + g(\epsilon, P). \quad (640)$$

After repeated use of Taylor's theorem we can collect all  $O(1)$ ,  $O(1/n)$  and  $O(1/\sqrt{n})$  terms into  $O(1)$ , and (634) follows. ■

*Theorem 67:* For the AWGN channel with SNR  $P$  and for  $0 < \epsilon \leq 1$ , we have

$$\log M_e^*(n, \epsilon, P) \geq nC - \sqrt{nV}Q^{-1}(\epsilon) + O(1). \quad (641)$$

Obtaining an expansion up to  $o(\sqrt{n})$  would only require Lemma 43. However, to refine the term to  $O(1)$  requires a certain lower bound on  $\kappa_\tau^n$  uniform in  $\tau \in [0, \delta)$  because we need to set  $\tau_n = O(1/\sqrt{n})$  instead of  $\tau = O(1)$ .

*Proof of Theorem 67:* We will use all the notation of the proof of Theorem 65, but redefine

$$\alpha_n = 1 - \epsilon + \frac{2B(P)}{\sqrt{n}}. \quad (642)$$

Note that for  $n$  sufficiently large  $\alpha_n < 1$  and the definition of  $\zeta_n$  in (623) is meaningful.

As in (625) we conclude that

$$\mathbb{P} \left[ \sum_{i=1}^n S_i \leq \zeta_n \right] \geq \alpha_n - \frac{B(P)}{\sqrt{n}} \geq 1 - \epsilon + \frac{B(P)}{\sqrt{n}}. \quad (643)$$

In other words, we have proven that for

$$\log \gamma_n = nC(P) - \zeta_n = nC(P) + \sqrt{nV(P)}Q^{-1}(\alpha_n) \quad (644)$$

we obtain

$$\begin{aligned} P_{Y^n|X=x^n} [i(x^n; Y^n) \geq \log \gamma_n] \\ = \mathbb{P} \left[ \sum_{i=1}^n S_i \leq \zeta_n \right] \geq 1 - \epsilon + \frac{B(P)}{\sqrt{n}} \end{aligned} \quad (645)$$

for sufficiently large  $n$  and any  $x^n \in F_n$ . Therefore, by setting

$$\tau_n \triangleq \frac{B(P)}{\sqrt{n}}. \quad (646)$$

we have

$$\begin{aligned} \log \beta_{1-\epsilon+\tau_n}^n \\ \leq P_{Y^n} [i(x^n; Y^n) \geq \log \gamma_n] \end{aligned} \quad (647)$$

$$= \mathbb{E} \left[ \exp\{-i(x^n; Y^n)\} \mathbf{1}_{\{i(x^n; Y^n) \geq \log \gamma_n\}} \mid X^n = x^n \right] \quad (648)$$

$$\leq -\log \gamma_n - \frac{1}{2} \log n + O(1) \quad (649)$$

$$= -nC(P) + \zeta_n - \frac{1}{2} \log n + O(1) \quad (650)$$

where the (649) is by Lemma 47.

Finally, we use general Theorem 25 with  $\tau = \tau_n$  to obtain

$$\log M_e^*(n, P, \epsilon) \geq \log \frac{\kappa_{\tau_n}^n}{\beta_{1-\epsilon+\tau_n}^n}. \quad (651)$$

For the chosen  $\tau_n$  Lemma 61 gives

$$\log \kappa_{\tau_n}^n \geq -\frac{1}{2} \log n + O(1). \quad (652)$$

This inequality, together with (650), yields

$$\log M_e^*(n, P, \epsilon) \geq nC(P) - \zeta_n + O(1). \quad (653)$$

It is easy to see that  $Q^{-1}(\alpha_n) = Q^{-1}(1 - \epsilon) + O(1/\sqrt{n})$  and, thus, for  $\zeta_n$  we have

$$\zeta_n = \sqrt{nV(P)}Q^{-1}(\epsilon) + O(1). \quad (654)$$

■

*Proof of Theorem 54:* Expansion (291) is implied by (294) and (295). The lower bounds in (294) and (295) follow from (641). The upper bound in (294) is given by (612) for equal-power constraint and by (633) for maximal-power constraint. The upper bound in (295) is proved by (634). ■

#### ACKNOWLEDGMENT

The authors are grateful to Dr. T. Richardson for generating the multiedge LDPC performance included in Fig. 12, and to Dr. A. Ashikhmin for suggesting the inclusion of the BSC bound in Theorem 4 and the BEC bound in Theorem 6.

#### REFERENCES

- [1] S. Verdú, "teaching it," in *Proc. XXVIII Shannon Lecture, 2007 IEEE Int. Symp. Inf. Theory*, Nice, France, Jun. 28, 2007.
- [2] S. Verdú and T. S. Han, "A general formula for channel capacity," *IEEE Trans. Inf. Theory*, vol. 40, pp. 1147–1157, 1994.
- [3] C. E. Shannon, "Probability of error for optimal codes in a Gaussian channel," *Bell Syst. Tech. J.*, vol. 38, pp. 611–656, 1959.
- [4] D. Slepian, "Bounds on communication," *Bell Syst. Tech. J.*, vol. 42, pp. 681–707, 1963.
- [5] S. Dolinar, D. Divsalar, and F. Pollara, Code Performance as a Function of Block Size Jet Propulsion Lab., Pasadena, CA, 1998, JPL TDA Progress Report, 42–133.
- [6] C. Salema, *Microwave Radio Links: From Theory to Design*. New York: Wiley, 2002.
- [7] D. Baron, M. A. Khojastepour, and R. G. Baraniuk, "How quickly can we approach channel capacity?," in *Proc. 38th Asilomar Conf. Signals, Syst., Comput.*, Pacific Grove, CA, Nov. 2004.
- [8] A. Valembois and M. P. C. Fossorier, "Sphere-packing bounds revisited for moderate block lengths," *IEEE Trans. Inf. Theory*, vol. 50, pp. 2998–3014, 2004.
- [9] D. E. Luzzi, T. Beth, and S. Egner, "Constrained capacity of the AWGN channel," in *Proc. 1998 IEEE Int. Symp. Inf. Theory (ISIT)*, Cambridge, MA, 1998.
- [10] J. Shi and R. D. Wesel, "A study on universal codes with finite block lengths," *IEEE Trans. Inf. Theory*, vol. 53, pp. 3066–3074, 2007.
- [11] G. Wiechman and I. Sason, "An improved sphere-packing bound for finite-length codes over symmetric memoryless channels," *IEEE Trans. Inf. Theory*, vol. 54, pp. 1962–1990, 2008.
- [12] A. E. Ashikhmin, A. Barg, and S. N. Litsyn, "A new upper bound on the reliability function of the Gaussian channel," *IEEE Trans. Inf. Theory*, vol. 46, pp. 1945–1961, 2000.

- [13] A. Feinstein, "A new basic theorem of information theory," *IRE Trans. Inf. Theory*, vol. 4, no. 4, pp. 2–22, 1954.
- [14] C. E. Shannon, "Certain results in coding theory for noisy channels," *Inf. Contr.*, vol. 1, pp. 6–25, 1957.
- [15] R. G. Gallager, "A simple derivation of the coding theorem and some applications," *IEEE Trans. Inf. Theory*, vol. 11, pp. 3–18, 1965.
- [16] R. G. Gallager, *Information Theory and Reliable Communication*. New York: Wiley, 1968.
- [17] G. Polytyrev, "Bounds on the decoding error probability of binary linear codes via their spectra," *IEEE Trans. Inf. Theory*, vol. 40, pp. 1284–1292, 1994.
- [18] A. Barg and G. D. Forney, "Random codes: Minimum distances and error exponents," *IEEE Trans. Inf. Theory*, vol. 48, pp. 2568–2573, 2002.
- [19] T. Helleseth, T. Klove, and V. Levenshtein, "On the information function of an error correcting code," *IEEE Trans. Inf. Theory*, vol. 43, pp. 549–557, 1997.
- [20] A. Ashikhmin, personal communication 2009.
- [21] A. J. Thomasian, "Error bounds for continuous channels," in *Proc. 4th London Symp. Inf. Theory*, C. Cherry, Ed., Washington, DC, 1961, pp. 46–60.
- [22] R. Ash, *Information Theory*. New York: Interscience, 1965.
- [23] J. Wolfowitz, "The coding of messages subject to chance errors," *Illinois J. Math.*, vol. 1, pp. 591–606, 1957.
- [24] J. Wolfowitz, *Coding Theorems of Information Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1962.
- [25] J. Wolfowitz, "Notes on a general strong converse," *Inf. Contr.*, vol. 12, pp. 1–4, 1968.
- [26] H. V. Poor and S. Verdú, "A lower bound on the error probability in multihypothesis testing," *IEEE Trans. Inf. Theory*, vol. 41, pp. 1992–1993, 1995.
- [27] C. E. Shannon, R. G. Gallager, and E. R. Berlekamp, "Lower bounds to error probability for coding on discrete memoryless channels I," *Inf. Contr.*, vol. 10, pp. 65–103, 1967.
- [28] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, pp. 379–423, Oct. 1948.
- [29] L. Weiss, "On the strong converse of the coding theorem for symmetric channels without memory," *Quart. Appl. Math.*, vol. 18–3, 1960.
- [30] R. L. Dobrushin, "Mathematical problems in the Shannon theory of optimal coding of information," in *Proc. 4th Berkeley Symp. Math., Statist., Probabil.*, 1961, vol. 1, pp. 211–252.
- [31] V. Strassen, "Asymptotische abschätzungen in Shannon's informationstheorie," in *Trans. 3rd Prague Conf. Inf. Theory*, Prague, 1962, pp. 689–723.
- [32] L. Wang, R. Colbeck, and R. Renner, "Simple channel coding bounds," in *Proc. 2009 IEEE Int. Symp. Inf. Theory (ISIT)*, Seoul, Korea, Jul. 2009.
- [33] Y. Polyanskiy, H. V. Poor, and S. Verdú, "New channel coding achievability bounds," in *Proc. 2008 IEEE Int. Symp. Inf. Theory (ISIT)*, Toronto, ON, Canada, 2008.
- [34] I. Csiszár, "Information-type measures of difference of probability distributions and indirect observation," *Studia Sci. Math. Hungar.*, vol. 2, pp. 229–318, 1967.
- [35] S. Verdú, *EE528—Information Theory, Lecture Notes*. Princeton, NJ: Princeton Univ., 2007.
- [36] R. E. Blahut, "Hypothesis testing and information theory," *IEEE Trans. Inf. Theory*, vol. 20, pp. 405–417, 1974.
- [37] S. J. MacMullan and O. M. Collins, "A comparison of known codes, random codes and the best codes," *IEEE Trans. Inf. Theory*, vol. 44, pp. 3009–3022, 1998.
- [38] C. Di, D. Proietti, I. E. Telatar, T. J. Richardson, and R. Urbanke, "Finite-length analysis of low-density parity-check codes on the binary erasure channel," *IEEE Trans. Inf. Theory*, vol. 48, pp. 1570–1579, 2002.
- [39] P. Elias, "Coding for two noisy channels," in *Proc. 3rd London Symp. Inf. Theory*, Washington, DC, Sep. 1955, pp. 61–76.
- [40] W. Feller, *An Introduction to Probability Theory and Its Applications*, Second ed. New York: Wiley, 1971, vol. II.
- [41] P. Van Beeck, "An application of Fourier methods to the problem of sharpening the Berry-Esseen inequality," *Z. Wahrscheinlichkeitstheorie und Verw. Geb.*, vol. 23, pp. 187–196, 1972.
- [42] H. V. Poor, *An Introduction to Signal Detection and Estimation*, 2nd ed. New York: Springer-Verlag, 1994.
- [43] M. Abramowitz and I. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, 10th ed. New York: Dover, 1972.
- [44] A. L. Jones, "An extension of inequality involving modified Bessel functions," *J. Math. Phys.*, vol. 47, pp. 220–221, 1968.
- [45] A. V. Prokhorov, "Inequalities for Bessel functions of a purely imaginary argument," *Theor. Probability Appl.*, vol. 13, pp. 496–501, 1968.
- [46] C.-G. Esseen, "On the concentration function of a sum of independent random variables," *Z. Wahrscheinlichkeitstheorie und Verw. Geb.*, vol. 9, no. 4, pp. 290–308, 1968.
- [47] M. Hayashi, "Information spectrum approach to second-order coding rate in channel coding," *IEEE Trans. Inf. Theory*, vol. 55, pp. 4947–4966, Nov. 2009.
- [48] D. Buckingham and M. C. Valenti, "The information-outage probability of finite-length codes over AWGN channels," in *Proc. Conf. Inf. Sci. Syst. (CISS)*, Princeton, NJ, Mar. 2008.
- [49] T. Richardson, personal communication 2009.
- [50] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Dispersion of the Gilbert-Elliott channel," in *Proc. 2009 IEEE Int. Symp. Inf. Theory (ISIT)*, Seoul, Korea, Jul. 2009.
- [51] C. E. Shannon, "Behavior near channel capacity," in *Claude Elwood Shannon's Papers*. Washington, DC: Manuscript Division, Library of Congress, unpublished.

**Yury Polyanskiy** (S'08) received the B.S. and M.S. degrees (both with honors) in applied mathematics and physics from the Moscow Institute of Physics and Technology in 2003 and 2005, respectively. He is currently pursuing a Ph.D. degree in electrical engineering at Princeton University, Princeton, NJ.

In 2000–2005, he was with the Department of Surface Oilfield Equipment, Borets Company LLC, where he rose to the position of Chief Software Designer. His research interests include information theory, coding theory and the theory of random processes.

Mr. Polyanskiy was a recipient of the Best Student Paper Award at the 2008 IEEE International Symposium on Information Theory (ISIT), Toronto, ON, Canada.

**H. Vincent Poor** (S'72–M'77–SM'82–F'87) received the Ph.D. degree in electrical engineering and computer science from Princeton University, Princeton, NJ, in 1977.

From 1977 until 1990, he was on the faculty of the University of Illinois at Urbana-Champaign. Since 1990, he has been on the faculty at Princeton University, where he is the Dean of Engineering and Applied Science, and the Michael Henry Strater University Professor of Electrical Engineering. His research interests are in the areas of stochastic analysis, statistical signal processing, and information theory, and their applications in wireless networks and related fields. Among his publications in these areas are the recent books *Quickest Detection* (Cambridge University Press, 2009), coauthored with O. Hadjilias and *Information Theoretic Security* (Now Publishers, 2009), coauthored with Y. Liang and S. Shamai.

Dr. Poor is a member of the National Academy of Engineering, a Fellow of the American Academy of Arts and Sciences, and an International Fellow of the Royal Academy of Engineering of the U.K. He is also a Fellow of the Institute of Mathematical Statistics, the Optical Society of America, and other organizations. In 1990, he served as President of the IEEE Information Theory Society, and in 2004–2007 as the Editor-in-Chief of these TRANSACTIONS. He was the recipient of the 2005 IEEE Education Medal. Recent recognition of his work includes the 2007 Technical Achievement Award of the IEEE Signal Processing Society, the 2008 Aaron D. Wyner Distinguished Service Award of the IEEE Information Theory Society, and the 2009 Edwin Howard Armstrong Achievement Award of the IEEE Communications Society.

**Sergio Verdú** (S'80–M'84–SM'88–F'93) received the Telecommunications Engineering degree from the Universitat Politècnica de Barcelona, Barcelona, Spain, in 1980 and the Ph.D. degree in Electrical Engineering from the University of Illinois at Urbana-Champaign, Urbana, in 1984.

Since 1984, he has been a member of the faculty of Princeton University, Princeton, NJ, where he is the Eugene Higgins Professor of Electrical Engineering.

Dr. Verdú is the recipient of the 2007 Claude E. Shannon Award and the 2008 IEEE Richard W. Hamming Medal. He is a member of the National Academy of Engineering and was awarded a Doctorate Honoris Causa from the Universitat Politècnica de Catalunya in 2005. He is a recipient of several paper awards from the IEEE: the 1992 Donald Fink Paper Award, the 1998 Information Theory Outstanding Paper Award, an Information Theory Golden Jubilee Paper Award, the 2002 Leonard Abraham Prize Award, the 2006 Joint Communications/Information Theory Paper Award, and the 2009 Stephen O. Rice Prize from IEEE Communications Society. He has also received paper awards from the Japanese Telecommunications Advancement Foundation and from Eurasip. He received the 2000 Frederick E. Terman Award from the American Society for Engineering Education for his book *Multiuser Detection* (Cambridge, U.K.: Cambridge Univ. Press, 1998). He served as President of the IEEE Information Theory Society in 1997. He is currently Editor-in-Chief of *Foundations and Trends in Communications and Information Theory*.