

# Exact Random Coding Error Exponents of Optimal Bin Index Decoding \*

Neri Merhav

Department of Electrical Engineering  
Technion - Israel Institute of Technology  
Technion City, Haifa 32000, ISRAEL  
E-mail: merhav@ee.technion.ac.il

## Abstract

We consider ensembles of channel codes that are partitioned into bins, and focus on analysis of exact random coding error exponents associated with optimum decoding of the index of the bin to which the transmitted codeword belongs. Two main conclusions arise from this analysis: (i) for independent random selection of codewords within a given type class, the random coding exponent of optimal bin index decoding is given by the ordinary random coding exponent function, computed at the rate of the entire code, independently of the exponential rate of the size of the bin. (ii) for this ensemble of codes, sub-optimal bin index decoding, that is based on ordinary maximum likelihood (ML) decoding, is as good as the optimal bin index decoding in terms of the random coding error exponent achieved. Finally, for the sake of completeness, we also outline how our analysis of exact random coding exponents extends to the hierarchical ensemble that correspond to superposition coding and optimal decoding, where for each bin, first, a cloud center is drawn at random, and then the codewords of this bin are drawn conditionally independently given the cloud center. For this ensemble, conclusions (i) and (ii), mentioned above, no longer hold necessarily in general.

**Index Terms:** Random coding, error exponent, binning, broadcast channels, superposition coding.

---

\*This research was supported by the Israel Science Foundation (ISF), grant no. 412/12.

# 1 Introduction

In multiuser information theory, one of the most frequently encountered building blocks is the notion of *superposition coding*, namely, coding with an hierarchical structure, in which the codebook is naturally partitioned into *bins*, or *clouds*. The original idea of superposition coding dates back to Cover [2], who proposed it in the context of broadcast channels (see also [3], [6, Section 15.6] and references therein). Later on, it has been proved extremely useful in a much wider variety of coded communication settings, including the wiretap channel [7], [17], the Gel'fand–Pinsker channel [9] (and in duality, Wyner–Ziv source encoding [18]), the relay channel [4], the interference channel [1], the multiple access channel [10], and channels with feedback [5], [12], just to name a few.

Generally speaking, the aim of superposition coding is to encode pairs of messages jointly, such that each message pair is mapped into a single codeword. To this end, the codebook is constructed with an hierarchical structure of bins (or clouds), such that a receiver that operates under relatively good channel conditions (high SNR) can decode reliably both messages, whereas a receiver that works under relatively bad channel conditions (low SNR) can decode reliably at least one of the messages, the one which consists of the index of the bin to which the codeword belongs.

This hierarchical structure of partitioning into bins is applicable even in achievability schemes with simple code ensembles, where all codewords are drawn independently under a certain distribution. Consider, for example, a random code of size  $M_1 = e^{nR_1}$ , where each codeword  $\mathbf{x}_i = (x_{1,i}, x_{2,i}, \dots, x_{n,i})$ ,  $i = 0, 1, \dots, M_1 - 1$ , is selected independently at random with a uniform distribution over a given type class. The code is then divided into  $M = e^{nR}$  ( $R \leq R_1$ ) bins  $\{\mathcal{C}_w\}_{w=0}^{M-1}$ ,  $\mathcal{C}_w = \{\mathbf{x}_{wM_2}, \mathbf{x}_{wM_2+1}, \dots, \mathbf{x}_{(w+1)M_2-1}\}$ , where  $M_2 = M_1/M = e^{n(R_1-R)} \triangleq e^{nR_2}$ . Assuming that the choice of the index  $i$  of the transmitted codeword is governed by the uniform distribution over  $\{0, 1, \dots, M_1 - 1\}$ , our focus, in this paper, will be on the user that decodes merely the index  $w$  of the bin  $\mathcal{C}_w$  that contains  $\mathbf{x}_i$ , namely,  $w = \lfloor i/M_2 \rfloor$ . This problem setting, including the above described random coding ensemble, is the very same as the one encountered from the viewpoint of the legitimate receiver in the achievability scheme of the wiretap channel model [17], as well as the decoder in the direct part of the Gel'fand–Pinsker channel [9].

Denoting the channel output vector by  $\mathbf{y} = (y_1, \dots, y_n)$  and the channel transition probability

function by  $P(\mathbf{y}|\mathbf{x})$ , the optimal bin index decoder is given by

$$w^*(\mathbf{y}) = \operatorname{argmax}_{0 \leq w \leq M-1} P(\mathbf{y}|\mathcal{C}_w) \quad (1)$$

where

$$P(\mathbf{y}|\mathcal{C}_w) \triangleq \frac{1}{M_2} \sum_{\mathbf{x} \in \mathcal{C}_w} P(\mathbf{y}|\mathbf{x}) = \frac{1}{M_2} \sum_{i=wM_2}^{(w+1)M_2-1} P(\mathbf{y}|\mathbf{x}_i). \quad (2)$$

Another, suboptimal decoder, which is natural to consider for bin index decoding, is the one that first estimates the index of the transmitted codeword using the ordinary maximum likelihood (ML) decoder, i.e.,  $\hat{i}_{\text{ML}}(\mathbf{y}) = \operatorname{argmax}_{0 \leq i \leq M_1-1} P(\mathbf{y}|\mathbf{x}_i)$ , and then decodes the bin index  $\hat{w}(\mathbf{y})$  as the one that includes that codeword, i.e.,

$$\hat{w}(\mathbf{y}) = \left\lfloor \frac{\hat{i}_{\text{ML}}(\mathbf{y})}{M_2} \right\rfloor. \quad (3)$$

In fact, the decoder of the achievability scheme of [17] is closer in spirit to (3) than to (1), except that the estimator of  $i$  can even be defined there in terms of joint typicality rather than in terms of ML decoding (in order to facilitate the analysis). According to the direct part of the coding theorem in [17], for memoryless channels, such a decoder is good enough (in spite of its sub-optimality) for achieving the maximum achievable information rate, just like decoder (1). It therefore seems conceivable that decoder (3) would achieve the same maximum rate too. Similar comments apply to the decoder of [9], as well as those in many other related works that involve superposition coding.

The question that we will address in this paper is what happens if we examine decoder (3), in comparison to decoder (1), under the more refined criterion of the error exponent as a function of the rates  $R_1$  and  $R_2$ . Would decoder (3) achieve the same optimal error exponent as the optimal decoder (1)?

By analyzing the exact random coding error exponent associated with decoder (1), in comparison to (3), for a given memoryless channel, we answer this question affirmatively, at least for the ensemble of codes described above, where each codeword is selected independently at random, under the uniform distribution within a given type class. In particular, our main result is that both decoders achieve the error exponent given by  $E_r(R_1)$ , independently of  $R_2$ , where  $E_r(\cdot)$  is the random coding error exponent function of ordinary ML decoding for the above defined ensemble. In other words, decoder (3) is essentially as good as the optimal decoder (1), not only from the viewpoint of achievable information rates, but moreover, in terms of error exponents.

The fact that the two decoders have the same error exponent may appear surprising at first glance. It indicates that for a considerable fraction<sup>1</sup> of the error events of (1), the score  $P(\mathbf{y}|\mathcal{C}_w)$  for a wrong bin may appear large (enough to exceed the one of the correct bin), mostly because of an incidental fluctuation in the likelihood of a single codeword (or a few codewords) within that bin, rather than due to a collective fluctuation of the entire bin (or a considerable fraction of it). Thus, it appears conceivable that many of the error events will be common to both decoders. Now, given that there is a single wrong codeword (in the entire codebook), whose likelihood is exceedingly large, the probability that it would belong to an incorrect bin is about  $(M - 1)/M = 1 - 1/M$  (due to symmetry), thus roughly speaking, erroneous bin index decoding is essentially as frequent as erroneous decoding of the ordinary ML decoder. Consequently, its probability depends on  $R_2$  so weakly that its asymptotic exponent is completely independent of  $R_2$ . This independence of  $R_2$  means that the reliability of decoding part of a message ( $nR$  out of  $nR_1$  nats) is essentially the same as that of decoding the entire message, no matter how small or large the size of this partial message may be.

The exponential equivalence of the performance of the two decoders should be interpreted as an encouraging message, because the optimal decoder (1) is extremely difficult to implement numerically, as the calculation of each score involves the summation of  $M_2$  terms  $\{P(\mathbf{y}|\mathbf{x}_i)\}$ , which are typically extremely small numbers for large  $n$  (usually obtained from long products of numbers between zero and one). On the other hand, decoder (3) easily lends itself to calculations in the logarithmic domain, where products are transformed into sums, thus avoiding these difficult numerical problems. Moreover, if the underlying memoryless channel is unknown, decoder (3) can easily be replaced by a similar decoder that is based on the universal maximum mutual information (MMI) decoder [8], while it is less clear how to transform (1) into a universal decoder.

Yet another advantage of decoder (3) is associated with the perspective of mismatch. Let the true underlying channel  $P(\mathbf{y}|\mathbf{x}_i)$  be replaced by an incorrect assumed channel  $P'(\mathbf{y}|\mathbf{x}_i)$ , both in (1) and (3). It turns out that the random coding error exponent of the latter is never worse (and sometimes may be better) than the former. Thus, decoder (3) is more robust to mismatch.

For the sake of completeness, we also extend our exact error exponent analysis to account for the hierarchical ensemble of superposition coding (applicable for the broadcast channel), where

---

<sup>1</sup>Namely, a fraction that maintains the exponential rate of the probability of the error event.

first,  $M$  cloud centers,  $\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{M-1}$ , are drawn independently at random from a given type class, and then for each  $\mathbf{u}_w$ ,  $w = 0, 1, \dots, M-1$ ,  $M_2$  codewords  $\mathbf{x}_{wM_2}, \mathbf{x}_{wM_2+1}, \dots, \mathbf{x}_{(w+1)M_2-1}$  are drawn conditionally independently from a given conditional type class given  $\mathbf{u}_w$ . The resulting error exponent is the exact<sup>2</sup> random coding exponent of the weak decoder in the degraded broadcast channel model. Here, it is no longer necessarily true that the error exponent is independent of  $R_2$  and that decoders (1) and (3) achieve the same exponent.

Finally, it should be pointed out that in a recent paper [14], a complementary study, of the random coding exponent of *correct* decoding, for the optimal bin index decoder, was carried out for rates above the maximum rate of reliable communication (i.e., the mutual information induced by the empirical distribution of the codewords and the channel). Thus, while this paper is relevant for the legitimate decoder of the wiretap channel model [17], the earlier work [14] is relevant for the decoder of the wiretapper of the same model.

The outline of the remaining part of this paper is as follows. In Section 2, we establish notation conventions. In Section 3, we formalize the problem and assert the main theorem concerning the error exponent of decoders (1) and (3). Section 4 is devoted to the proof of this theorem, and in Section 5, we discuss it. Finally, in Section 6, we extend our error exponent analysis to the case where the ensemble of random codes is defined hierarchically.

## 2 Notation Conventions

Throughout the paper, random variables will be denoted by capital letters, specific values they may take will be denoted by the corresponding lower case letters, and their alphabets will be denoted by calligraphic letters. Random vectors and their realizations will be denoted, respectively, by capital letters and the corresponding lower case letters, both in the bold face font. Their alphabets will be superscripted by their dimensions. For example, the random vector  $\mathbf{X} = (X_1, \dots, X_n)$ , ( $n$  – positive integer) may take a specific vector value  $\mathbf{x} = (x_1, \dots, x_n)$  in  $\mathcal{X}^n$ , the  $n$ -th order Cartesian power of  $\mathcal{X}$ , which is the alphabet of each component of this vector. The probability of an event  $\mathcal{E}$  will be denoted by  $\Pr\{\mathcal{E}\}$ , and the expectation operator will be denoted by  $\mathbf{E}\{\cdot\}$ . For two positive sequences  $a_n$  and  $b_n$ , the notation  $a_n \doteq b_n$  will stand for equality in the exponential scale, that

---

<sup>2</sup>This is different from earlier work (see [11] and references therein), where lower bounds were derived.

is,  $\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{a_n}{b_n} = 0$ . Thus,  $a_n \doteq 0$  means that  $a_n$  tends to zero in a super-exponential rate. Similarly,  $a_n \leq b_n$  means that  $\limsup_{n \rightarrow \infty} \frac{1}{n} \log \frac{a_n}{b_n} \leq 0$ , and so on. The indicator function of an event  $\mathcal{E}$  will be denoted by  $\mathcal{I}\{\mathcal{E}\}$ . The notation  $[x]_+$  will stand for  $\max\{0, x\}$ . Logarithms and exponents will be understood to be taken to the natural base unless specified otherwise.

Probability distributions, associated with sources and channels, will be denoted by the letters  $P$  and  $Q$ , with subscripts that denote the names of the random variables involved along with their conditioning, if applicable, following the customary notation rules in probability theory. For example,  $Q_{XY}$  stands for a generic joint distribution  $\{Q_{XY}(x, y), x \in \mathcal{X}, y \in \mathcal{Y}\}$ ,  $P_{Y|X}$  denotes the matrix of single-letter transition probabilities of the underlying memoryless channel from  $X$  to  $Y$ ,  $\{P_{Y|X}(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$ , and so on. Information measures induced by the generic joint distribution  $Q_{XY}$ , or  $Q$  for short, will be subscripted by  $Q$ , for example,  $I_Q(X; Y)$  will denote the corresponding mutual information, etc. The weighted divergence between two channels,  $Q_{Y|X}$  and  $P_{Y|X}$ , with weight  $P_X$ , is defined as

$$D(Q_{Y|X} \| P_{Y|X} | P_X) \triangleq \sum_{x \in \mathcal{X}} P_X(x) \sum_{y \in \mathcal{Y}} Q_{Y|X}(y|x) \ln \frac{Q_{Y|X}(y|x)}{P_{Y|X}(y|x)}. \quad (4)$$

The type class,  $\mathcal{T}(P_X)$ , associated with a given empirical probability distribution  $P_X$  of  $X$ , is the set of all  $\mathbf{x} = (x_1, \dots, x_n)$ , whose empirical distribution is  $P_X$ . Similarly, the joint type class of pairs of sequences  $\{(\mathbf{u}, \mathbf{x})\}$  in  $\mathcal{U}^n \times \mathcal{X}^n$ , which is associated with an empirical joint distribution  $P_{UX}$ , will be denoted by  $\mathcal{T}(P_{UX})$ . Finally, for a given  $P_{X|U}$  and  $\mathbf{u} \in \mathcal{U}^n$ ,  $\mathcal{T}(P_{X|U} | \mathbf{u})$  denotes the conditional type class of  $\mathbf{x}$  given  $\mathbf{u}$  w.r.t.  $P_{X|U}$ , namely, the set of sequences  $\{\mathbf{x}\}$  whose conditional empirical distribution w.r.t.  $\mathbf{u}$  is given by  $P_{X|U}$ .

### 3 Problem Formulation and Main Result

Consider a discrete memoryless channel (DMC), defined by a matrix of single-letter transition probabilities,  $\{P_{Y|X}(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are finite alphabets. When the channel is fed with an input vector  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$ , the output is a random vector  $\mathbf{Y} = (Y_1, \dots, Y_n) \in \mathcal{Y}^n$ , distributed according to

$$P(\mathbf{y} | \mathbf{x}) = \prod_{t=1}^n P_{Y|X}(y_t | x_t), \quad (5)$$

where to avoid cumbersome notation, here and throughout the sequel, we omit the subscript “ $\mathbf{Y}|\mathbf{X}$ ” in the notation of the conditional distribution of the vector channel, from  $\mathcal{X}^n$  to  $\mathcal{Y}^n$ . Consider next a codebook,  $\mathcal{C} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{M_1-1}\}$ , where  $M_1 = e^{nR_1}$ , and where each  $\mathbf{x}_i$ ,  $i = 0, 1, \dots, M_1 - 1$ , is selected independently at random, under the uniform distribution across the type class  $\mathcal{T}(P_X)$ , where  $P_X$  is a given distribution over  $\mathcal{X}$ . Once selected, the codebook  $\mathcal{C}$  is revealed to both the encoder and the decoder. The codebook  $\mathcal{C}$  is partitioned into  $M = e^{nR}$  bins,  $\{\mathcal{C}_w\}_{w=0}^{M-1}$ , each one of size  $M_2 = e^{nR_2}$  ( $R + R_2 = R_1$ ), where  $\mathcal{C}_w = \{\mathbf{x}_{wM_2}, \mathbf{x}_{wM_2+1}, \dots, \mathbf{x}_{(w+1)M_2-1}\}$ .

Let  $\mathbf{x}_I \in \mathcal{C}$  be transmitted over the channel, where  $I$  is a random variable drawn under the uniform distribution over  $\{0, 1, \dots, M_1 - 1\}$ , independently of the random selection of the code. Let  $W = \lfloor I/M_2 \rfloor$  designate the random bin index to which  $\mathbf{x}_I$  belongs and let  $\mathbf{Y} \in \mathcal{Y}^n$  be the channel output resulting from the transmission of  $\mathbf{x}_I$ .

Consider the bin index decoders (1) and (3), and define their average error probabilities, as

$$P_e^* = \mathbf{E}[\Pr\{w^*(\mathbf{Y}) \neq W\}], \quad \hat{P}_e = \mathbf{E}[\Pr\{\hat{w}(\mathbf{Y}) \neq W\}], \quad (6)$$

where the probabilities are defined w.r.t. the randomness of the index  $I$  of the transmitted codeword (hence the randomness of  $W$ ) and the random operation of the channel, and the expectations are taken w.r.t. the randomness of the codebook  $\mathcal{C}$ .

Our goal is to assess the exact exponential rates of  $P_e^*$  and  $\hat{P}_e$ , as functions of  $R_1$  and  $R_2$ , that is,

$$E^*(R_1, R_2) \triangleq \lim_{n \rightarrow \infty} \left[ -\frac{\ln P_e^*}{n} \right] \quad (7)$$

and

$$\hat{E}(R_1, R_2) \triangleq \lim_{n \rightarrow \infty} \left[ -\frac{\ln \hat{P}_e}{n} \right]. \quad (8)$$

At this point, a technical comment is in order. The case  $R = 0$  ( $R_1 = R_2$ ) should not be understood as a situation where there is only one bin and  $\mathcal{C}_0 = \mathcal{C}$ , since this is a degenerated situation, where there is nothing to decode as far as bin index decoding is concerned, the probability of error is trivially zero (just like in ordinary decoding, where there is only one codeword, which is meaningless). The case  $R = 0$  should be understood as a case where the number of bins is at least two, and at most sub-exponential in  $n$ . On the other extreme, for  $R_2 = 0$  ( $R_1 = R$ ), it is safe to consider each bin as consisting of a single codeword, rendering the case of ordinary decoding as a special case.

Our main result is the following.

**Theorem 1** *Let  $R_1$  and  $R_2$  be given ( $R_2 \leq R_1$ ). Let  $E^*(R_1, R_2)$  and  $\hat{E}(R_1, R_2)$  be defined as in eqs. (7) and (8), respectively. Then,*

$$E^*(R_1, R_2) = \hat{E}(R_1, R_2) = E_r(R_1) \quad (9)$$

where  $E_r(R_1)$  is the random coding error exponent function, i.e.,

$$E_r(R_1) = \min_{Q_{XY}: Q_X=P_X} \{D(Q_{Y|X} \| P_{Y|X} | P_X) + [I_Q(X; Y) - R_1]_+\}. \quad (10)$$

## 4 Proof of Theorem 1

For a given  $\mathbf{y} \in \mathcal{Y}^n$ , and a given joint probability distribution  $Q_{XY}$  on  $\mathcal{X} \times \mathcal{Y}$ , let  $N_1(Q_{XY})$  denote the number of codewords  $\{\mathbf{x}_i\}$  in  $\mathcal{C}_1$  whose conditional empirical joint distribution with  $\mathbf{y}$  is  $Q_{XY}$ , that is

$$N_1(Q_{XY}) = \sum_{i=M_2}^{2M_2-1} \mathcal{I}\{(\mathbf{x}_i, \mathbf{y}) \in \mathcal{T}(Q_{XY})\}. \quad (11)$$

We also denote

$$f(Q_{XY}) = \frac{1}{n} \ln P(\mathbf{y}|\mathbf{x}) = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} Q_{XY}(x, y) \ln P_{Y|X}(y|x), \quad (12)$$

where  $Q_{XY}$  is understood to be the joint empirical distribution of  $(\mathbf{x}, \mathbf{y}) \in \mathcal{X}^n \times \mathcal{Y}^n$ . Without loss of generality, we assume throughout, that the transmitted codeword is  $\mathbf{x}_0$ , and so, the correct bin is  $\mathcal{C}_0$ . The average probability of error, associated with decoder (1), is given by

$$P_e^* \doteq \mathbf{E} [\min\{1, M \cdot \Pr\{P(\mathbf{Y}|\mathcal{C}_1) \geq P(\mathbf{Y}|\mathcal{C}_0)\}\}], \quad (13)$$

where the expectation is w.r.t. the randomness of  $\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_{M_2-1}$  and  $\mathbf{Y}$ , and where given  $\mathbf{X}_0 = \mathbf{x}_0$ ,  $\mathbf{Y}$  is distributed according to  $P(\cdot|\mathbf{x}_0)$ . For a given  $\mathbf{y}$ , the pairwise error probability,  $\Pr\{P(\mathbf{y}|\mathcal{C}_1) \geq P(\mathbf{y}|\mathcal{C}_0)\}$ , is calculated w.r.t. the randomness of  $\mathcal{C}_1 = \{\mathbf{X}_{M_2}, \mathbf{X}_{M_2+1}, \dots, \mathbf{X}_{2M_2-1}\}$ , but for a given  $\mathcal{C}_0$ . To see why (13) is true, first observe that the right-hand side (r.h.s.) of this equation is simply the (expectation of the) union bound, truncated to unity. On the other hand, since the pairwise error events are conditionally independent given  $\mathcal{C}_0$  and  $\mathbf{y}$ , the r.h.s. times a



factor of 1/2 (which does not affect the exponent), serves as a lower bound to the probability of the union of the pairwise error events [15, Lemma A.2] (see also [16, Lemma 1]).

We next move on to the calculation of the pairwise error probability. For a given  $\mathcal{C}_0$  and  $\mathbf{y}$ , let

$$s \triangleq \frac{1}{n} \ln \left[ \sum_{i=0}^{M_2-1} P(\mathbf{y}|\mathbf{x}_i) \right], \quad (14)$$

and so, the pairwise error probability becomes  $\Pr\{M_2 \cdot P(\mathbf{y}|\mathcal{C}_1) \geq e^{ns}\}$ , where it is kept in mind that  $s$  is a function of  $\mathcal{C}_0$  and  $\mathbf{y}$ . Now,

$$\Pr\{M_2 \cdot P(\mathbf{y}|\mathcal{C}_1) \geq e^{ns}\} = \Pr \left\{ \sum_{i=M_2}^{2M_2-1} P(\mathbf{y}|\mathbf{X}_i) \geq e^{ns} \right\} \quad (15)$$

$$= \Pr \left\{ \sum_{Q_{X|Y}} N_1(Q_{XY}) e^{nf(Q_{XY})} \geq e^{ns} \right\} \quad (16)$$

$$\doteq \Pr \left\{ \max_{Q_{X|Y} \in \mathcal{S}(Q_Y)} N_1(Q_{XY}) e^{nf(Q_{XY})} \geq e^{ns} \right\} \quad (17)$$

$$= \Pr \bigcup_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \left\{ N_1(Q_{XY}) e^{nf(Q_{XY})} \geq e^{ns} \right\} \quad (18)$$

$$\doteq \sum_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \Pr \left\{ N_1(Q_{XY}) e^{nf(Q_{XY})} \geq e^{ns} \right\} \quad (19)$$

$$\doteq \max_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \Pr \left\{ N_1(Q_{XY}) \geq e^{n[s-f(Q_{XY})]} \right\}, \quad (20)$$

where for a given  $Q_Y$ ,  $\mathcal{S}(Q_Y)$  is defined as the set of all  $\{Q_{X|Y}\}$ , such that  $\sum_{y \in \mathcal{Y}} Q_Y(y) Q_{X|Y}(x|y) = P_X(x)$  for all  $x \in \mathcal{X}$ . Now, for a given  $Q_{XY}$ ,  $N_1(Q_{XY})$  is a binomial random variable with  $e^{nR_2}$  trials and probability of ‘success’ which is of the exponential order of  $e^{-nI_Q(X;Y)}$ . Thus, a standard large deviations analysis (see, e.g., [13, pp. 167–169]) yields

$$\Pr \left\{ N_1(Q_{X|Y}) \geq e^{n[s-f(Q_{XY})]} \right\} \doteq e^{-nE_0(Q_{XY})}, \quad (21)$$

where

$$E_0(Q_{XY}) = \begin{cases} [I(Q_{XY}) - R_2]_+ & f(Q_{XY}) \geq s \\ 0 & f(Q_{XY}) < s, f(Q_{XY}) \geq s - R_2 + I(Q_{XY}) \\ \infty & f(Q_{XY}) < s, f(Q_{XY}) < s - R_2 + I(Q_{XY}) \end{cases} \quad (22)$$

$$= \begin{cases} I(Q_{XY}) - R_2 & f(Q_{XY}) \geq s, R_2 < I(Q_{XY}) \\ 0 & f(Q_{XY}) \geq s, R_2 \geq I(Q_{XY}) \\ 0 & f(Q_{XY}) < s, f(Q_{XY}) \geq s - R_2 + I(Q_{XY}) \\ \infty & f(Q_{XY}) < s, f(Q_{XY}) < s - R_2 + I(Q_{XY}) \end{cases} \quad (23)$$

$$= \begin{cases} I(Q_{XY}) - R_2 & f(Q_{XY}) \geq s, R_2 < I(Q_{XY}) \\ 0 & f(Q_{XY}) \geq s - [R_2 - I(Q_{XY})]_+, R_2 \geq I(Q_{XY}) \\ \infty & f(Q_{XY}) < s - [R_2 - I(Q_{XY})]_+ \end{cases} \quad (24)$$

$$= \begin{cases} [I(Q_{XY}) - R_2]_+ & f(Q_{XY}) \geq s - [R_2 - I(Q_{XY})]_+ \\ \infty & f(Q_{XY}) < s - [R_2 - I(Q_{XY})]_+ \end{cases} \quad (25)$$

Therefore,  $\max_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \Pr\{N_1(Q_{XY}) \geq e^{n[s-f(Q_{XY})]}\}$  decays according to

$$E_1(s, Q_Y) = \min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} E_0(Q_{XY}),$$

which is given by

$$E_1(s, Q_Y) = \min\{[I_Q(X; Y) - R_2]_+ : f(Q_{XY}) + [R_2 - I_Q(X; Y)]_+ \geq s\} \quad (26)$$

with the understanding that the minimum over an empty set is defined as infinity. Finally,

$$P_e^* \doteq \mathbf{E} \min\{1, M \cdot e^{-nE_1(S, Q_Y)}\} = \mathbf{E} e^{-n[E_1(S, Q_Y) - R]_+}, \quad (27)$$

where the expectation is w.r.t. to the randomness of

$$S = \frac{1}{n} \ln \left[ \sum_{i=0}^{M_2-1} P(\mathbf{Y}|\mathbf{X}_i) \right] \quad (28)$$

and the randomness of  $Q_Y$ , the empirical distribution of  $\mathbf{Y}$ . This expectation will be taken in two steps: first, over the randomness of  $\{\mathbf{X}_1, \dots, \mathbf{X}_{M_2-1}\}$  while  $\mathbf{X}_0 = \mathbf{x}_0$  (the real transmitted codeword) and  $\mathbf{Y} = \mathbf{y}$  are held fixed, and then over the randomness of  $\mathbf{X}_0$  and  $\mathbf{Y}$ . Let  $\mathbf{x}_0$  and  $\mathbf{y}$  be given and let  $\epsilon > 0$  be arbitrarily small. Then,

$$\begin{aligned} P_e(\mathbf{x}_0, \mathbf{y}_0) &\triangleq \mathbf{E}[\exp\{-n[E_1(S, Q_Y) - R]_+\} | \mathbf{X}_0 = \mathbf{x}_0, \mathbf{Y} = \mathbf{y}] \\ &= \sum_s \Pr\{S = s | \mathbf{X}_0 = \mathbf{x}_0, \mathbf{Y} = \mathbf{y}\} \cdot \exp\{-n[E_1(s, Q_Y) - R]_+\} \\ &\leq \sum_i \Pr\{i\epsilon \leq S < (i+1)\epsilon | \mathbf{X}_0 = \mathbf{x}_0, \mathbf{Y} = \mathbf{y}\} \cdot \exp\{-n[E_1(i\epsilon, Q_Y) - R]_+\}, \end{aligned} \quad (29)$$

where  $i$  ranges from  $\frac{1}{n\epsilon} \ln P(\mathbf{y}|\mathbf{x}_0)$  to  $R_2/\epsilon$ . Now,

$$\begin{aligned} e^{ns} &= P(\mathbf{y}|\mathbf{x}_0) + \sum_{i=1}^{M_2-1} P(\mathbf{y}|\mathbf{X}_i) \\ &= e^{nf(Q_{X_0Y})} + \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})}, \end{aligned} \quad (30)$$

where  $Q_{X_0Y}$  is the empirical distribution of  $(\mathbf{x}_0, \mathbf{y})$  and  $N_0(Q_{XY})$  is the number of codewords in  $\mathcal{C}_0 \setminus \{\mathbf{x}_0\}$  whose joint empirical distribution with  $\mathbf{y}$  is  $Q_{XY}$ . The first term in the second line of (30) is fixed at this stage. As for the second term, we have (similarly as before):

$$\Pr \left\{ \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})} \geq e^{nt} \right\} \doteq e^{-nE_1(t, Q_Y)}. \quad (31)$$

On the other hand,

$$\Pr \left\{ \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})} \leq e^{nt} \right\} \doteq \Pr \bigcap_{Q_{X|Y}} \left\{ N_0(Q_{XY}) \leq e^{n[t-f(Q_{XY})]} \right\}. \quad (32)$$

Now, if there exists at least one  $Q_{X|Y} \in \mathcal{S}(Q_Y)$  for which  $I_Q(X; Y) < R_2$  and  $R_2 - I_Q(X; Y) > t - f(Q_{XY})$ , then this  $Q_{X|Y}$  alone is responsible for a double exponential decay of the probability of the event  $\{N_0(Q_{XY}) \leq e^{n[t-f(Q_{XY})]}\}$ , let alone the intersection over all  $Q_{X|Y} \in \mathcal{S}(Q_Y)$ . On the other hand, if for every  $Q_{X|Y} \in \mathcal{S}(Q_Y)$ , either  $I_Q(X; Y) \geq R_2$  or  $R_2 - I_Q(X; Y) \leq t - f(Q_{XY})$ , then we have an intersection of polynomially many events whose probabilities all tend to unity. Thus, the probability in question behaves exponentially like an indicator function of the condition that for every  $Q_{X|Y} \in \mathcal{S}(Q_Y)$ , either  $I_Q(X; Y) \geq R_2$  or  $R_2 - I_Q(X; Y) \leq t - f(Q_{XY})$ , or equivalently,

$$\Pr \left\{ \sum_{Q_{XY}} N_0(Q_{XY}) e^{nf(Q_{XY})} \leq e^{nt} \right\} \doteq \mathcal{I} \left\{ R_2 \leq \min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \{I_Q(X; Y) + [t - f(Q_{XY})]_+\} \right\}. \quad (33)$$

Let us now find what is the minimum value of  $t$  for which the value of this indicator function is unity. The condition is equivalent to

$$\min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \max_{0 \leq a \leq 1} \{I_Q(X; Y) + a[t - f(Q_{XY})]\} \geq R_2 \quad (34)$$

or, equivalently:

$$\forall Q_{X|Y} \in \mathcal{S}(Q_Y) \exists 0 \leq a \leq 1 : I_Q(X; Y) + a[t - f(Q_{XY})] \geq R_2, \quad (35)$$

which can also be written as

$$\forall Q_{X|Y} \in \mathcal{S}(Q_Y) \exists 0 \leq a \leq 1 : t \geq f(Q_{XY}) + \frac{R_2 - I_Q(X; Y)}{a} \quad (36)$$

or equivalently,

$$t \geq \max_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \min_{0 \leq a \leq 1} \left[ f(Q_{XY}) + \frac{R_2 - I_Q(X; Y)}{a} \right] \quad (37)$$

$$= \max_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \left[ f(Q_{XY}) + \begin{cases} R_2 - I_Q(X; Y) & R_2 \geq I_Q(X; Y) \\ -\infty & R_2 < I_Q(X; Y) \end{cases} \right] \quad (38)$$

$$= R_2 + \max_{\{Q_{X|Y} \in \mathcal{S}(Q_Y) : I_Q(X; Y) \leq R_2\}} [f(Q_{XY}) - I_Q(X; Y)] \quad (39)$$

$$\triangleq s_0(Q_Y). \quad (40)$$

Similarly, it is easy to check that  $E_1(t, Q_Y)$  vanishes for  $t \leq s_0(Q_Y)$ . Thus, in summary, we have

$$\Pr \left\{ e^{nt} \leq \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})} \leq e^{n(t+\epsilon)} \right\} \doteq \begin{cases} 0 & t < s_0(Q_Y) - \epsilon \\ e^{-nE(t, Q_Y)} & t \geq s_0(Q_Y) \end{cases} \quad (41)$$

Therefore, for a given  $(\mathbf{x}_0, \mathbf{y})$ , the expected error probability w.r.t.  $\{\mathbf{X}_1, \dots, \mathbf{X}_{M_2-1}\}$  yields

$$P_e(\mathbf{x}_0, \mathbf{y}) = \mathbf{E} \{ e^{-n[E_1(S, Q_Y) - R]_+} | \mathbf{X}_0 = \mathbf{x}_0, \mathbf{Y} = \mathbf{y} \} \quad (42)$$

$$\leq \sum_i \Pr \left\{ e^{ni\epsilon} \leq \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})} \leq e^{n(i+1)\epsilon} \right\} \times \exp(-n[E_1(\max\{i\epsilon, f(Q_{X_0Y})\}, Q_Y) - R]_+) \quad (43)$$

$$\leq \sum_{i \geq s_0(Q_Y)/\epsilon} \exp\{-nE_1(i\epsilon, Q_Y)\} \cdot \exp(-n[E_1(\max\{i\epsilon, f(Q_{X_0Y})\}, Q_Y) - R]_+), \quad (44)$$

where the expression  $\max\{i\epsilon, f(Q_{X_0Y})\}$  in the argument of  $E_1(\cdot, Q_Y)$  is due to the fact that

$$S = \frac{1}{n} \ln \left[ e^{nf(Q_{X_0Y})} + \sum_{Q_{X|Y}} N_0(Q_{XY}) e^{nf(Q_{XY})} \right] \quad (45)$$

$$\geq \frac{1}{n} \ln [e^{nf(Q_{X_0Y})} + e^{ni\epsilon}] \quad (46)$$

$$\doteq \max\{i\epsilon, f(Q_{X_0Y})\}. \quad (47)$$

By using the fact that  $\epsilon$  is arbitrarily small, we obtain

$$P_e(\mathbf{x}_0, \mathbf{y}) \doteq \exp(-n[E_1(\max\{s_0(Q_Y), f(Q_{X_0Y})\}, Q_Y) - R]_+), \quad (48)$$

since the dominant contribution to the sum over  $i$  is due to the term  $i = s_0(Q_Y)/\epsilon$  (by the non-increasing monotonicity of the function  $E_1(\cdot, Q_Y)$ ). Denoting  $s_1(Q_{X_0Y}) = \max\{s_0(Q_Y), f(Q_{X_0Y})\}$ , we then have, after averaging w.r.t.  $(\mathbf{X}_0, \mathbf{Y})$ ,

$$E^*(R_1, R_2) = \min_{Q_{Y|X_0}} \{ D(Q_{Y|X_0} \| P_{Y|X_0} | P_{X_0}) + [E_1(s_1(Q_{X_0Y}), Q_Y) - R]_+ \}, \quad (49)$$

where the random variable  $X_0$  is a replica of  $X$ , that is,  $P_{X_0} = P_X$ .

We next simplify the formula of  $E^*(R_1, R_2)$ . Clearly,

$$E(s_1(Q_{X_0Y}), Q_Y) = E_1(\max\{s_0(Q_Y), f(Q_{X_0Y})\}, Q_Y) \quad (50)$$

$$= \max\{E_1(s_0(Q_Y), Q_Y), E_1(f(Q_{X_0Y}), Q_Y)\} \quad (51)$$

$$= \max\{0, E_1(f(Q_{X_0Y}), Q_Y)\} \quad (52)$$

$$= E_1(f(Q_{X_0Y}), Q_Y). \quad (53)$$

Therefore,

$$E^*(R_1, R_2) = \min_{Q_{Y|X_0}} \{D(Q_{Y|X_0} \| P_{Y|X_0} | P_X) + [E_1(f(Q_{X_0Y}), Q_Y) - R]_+\}. \quad (54)$$

Finally, using the simple identity  $[[x - a]_+ - b]_+ \equiv [x - a - b]_+$  ( $b \geq 0$ ), we can slightly simplify this expression to be

$$E^*(R_1, R_2) = \min_{Q_{Y|X_0}} \{D(Q_{Y|X_0} \| P_{Y|X_0} | P_X) + [I_0(Q_{X_0Y}) - R_1]_+\}, \quad (55)$$

where

$$I_0(Q_{X_0Y}) \triangleq \min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \{I_Q(X; Y) : f(Q_{XY}) + [R_2 - I_Q(X; Y)]_+ \geq f(Q_{X_0Y})\}. \quad (56)$$

Now, let us define

$$E'_r(R_1) \triangleq \min_{Q_{Y|X_0}} \{D(Q_{Y|X_0} \| P_{Y|X_0} | P_X) + [I'_0(Q_{X_0Y}) - R_1]_+\}, \quad (57)$$

where

$$I'_0(Q_{X_0Y}) = \min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \{I_Q(X; Y) : f(Q_{XY}) \geq f(Q_{X_0Y})\}. \quad (58)$$

At this point,  $E'_r(R_1)$  is readily identified as the ordinary random coding error exponent associated with ML decoding (i.e., the special case of  $E^*(R_1, R_2)$  where  $R_2 = 0$ ), which is known [8, p. 165, Theorem 5.2] to be identical to the random coding error exponent,  $E_r(R_1)$ , achieved by maximum mutual information (MMI) universal decoding, defined similarly, except that  $I'_0(Q_{X_0Y})$  is replaced by

$$I''_0(Q_{X_0Y}) = \min_{Q_{X|Y} \in \mathcal{S}(Q_Y)} \{I_Q(X; Y) : I_Q(X; Y) \geq I_Q(X_0; Y)\} = I_Q(X_0; Y), \quad (59)$$

thus leading to equivalence with eq. (10).

To complete the proof, we now argue that  $E_r(R_1) = E^*(R_1, R_2) = \hat{E}(R_1, R_2)$ . The inequality  $E_r(R_1) \equiv E'_r(R_1) \geq E^*(R_1, R_2)$  is obvious since the minimization that defines  $I'_0(Q_{X_0Y})$  is over a

smaller set of distributions than the one that defines  $I_0(Q_{X_0Y})$ . On the other hand, the converse inequality,  $E_r(R_1) \leq E^*(R_1, R_2)$ , is also true because of the following consideration. We claim that

$$E_r(R_1) \leq E'(R_1, R_2) \leq \hat{E}(R_1, R_2) \leq E^*(R_1, R_2), \quad (60)$$

where definitions and explanations are now in order: As defined,  $E_r(R_1)$  is the random coding error exponent associated with ordinary ML decoding and the ordinary probability of error for a random code at rate  $R_1$ . Now, let  $E'(R_1, R_2)$  be defined as the random coding exponent of the ML decoder, where only errors associated with winning codewords that are outside the correct bin  $\mathcal{C}_0$  are counted. In other words, assuming that  $\mathbf{x}_0$  was transmitted, this is the exponent of the probability of the event  $\{\max_{i \geq M_2} P(\mathbf{y}|\mathbf{x}_i) \geq P(\mathbf{y}|\mathbf{x}_0)\}$ . Since this error event is a subset of the ordinary error event, its exponent is at least as large as  $E_r(R_1)$ , hence the first inequality. Now,  $\hat{E}(R_1, R_2)$ , which is the error exponent of decoder (3), is in fact, the exponent of the probability of the event  $\{\max_{i \geq M_2} P(\mathbf{y}|\mathbf{x}_i) \geq \max_{i < M_2} P(\mathbf{y}|\mathbf{x}_i)\}$  (given  $\mathbf{x}_0$ ), which in turn is a subset of the previous error event defined, hence the second inequality. Finally, the last inequality follows from the optimality of decoder (1), whose error exponent cannot be smaller than that of (3). Thus, we conclude that all inequalities are, in fact, equalities, and so,  $E^*(R_1, R_2) = \hat{E}(R_1, R_2) = E_r(R_1)$ , completing the proof of Theorem 1.

## 5 Discussion

When  $R_2 = 0$ , that is, a subexponential number of codewords within each bin, Theorem 1 is actually not surprising since  $\sum_{\mathbf{x} \in \mathcal{C}_w} P(\mathbf{y}|\mathbf{x}) \doteq \max_{\mathbf{x} \in \mathcal{C}_w} P(\mathbf{y}|\mathbf{x})$ , but for  $R_2 > 0$ , the results are not quite trivial (at least for the author of this article). As explained in the Introduction, the intuition is that the error probability is dominated by few codewords within some bin, whose likelihood score is exceptionally high. Note also that bin index decoding is different from the situation in list decoding, where even for a subexponential list size, the error exponent is improved. This is not surprising, because in list decoding, the list depends on the likelihood scores, and they are not given by a fixed bin, which is arbitrary.

Theorem 1 tells us that, under the ordinary random coding regime, decoding only a part of a message (say, a header of  $nR$  nats out of the total of  $nR_1$ ) is as reliable as decoding the entire message, as far as error exponents go. As discussed in the Introduction, decoder (3) is easier to

implement. It is also clear how to universalize this decoder: for an unknown DMC, replace  $\hat{i}_{\text{ML}}(\mathbf{y})$  in (3) by

$$\hat{i}_{\text{MMI}}(\mathbf{y}) = \operatorname{argmax}_{0 \leq i \leq M_1 - 1} I_Q(X_i; Y), \quad (61)$$

where  $I_Q(X_i; Y)$  designates the empirical mutual information induced by  $(\mathbf{x}_i, \mathbf{y})$ . This universal bin index decoder still achieves  $E_r(R_1)$ .

As for the mismatched case, the only change in the derivation in Section 4 is that the definition of the function  $f(Q_{XY})$  is changed to  $f(Q_{XY}) = \sum_{x,y} Q(x, y) \ln P'_{Y|X}(y|x)$  (or more generally, to an arbitrary function of  $Q_{XY}$ ), where  $P'_{Y|X}(y|x)$  is the mismatched channel. Here, it will still be true that  $E'_r(R_1)$  defined as in (57) (but with  $f$  being redefined) is not smaller than the corresponding  $E^*(R_1, R_2)$ , but the converse inequality (that was leading to equality before) can no longer be claimed since it was based on the optimality of decoder (1), but now both decoders are suboptimal. This means that, for the purpose of bin index decoding, decoder (3), but with  $P_{Y|X}$  replaced by  $P'_{Y|X}$ , is never worse than the corresponding mismatched version of decoder (1).

## 6 Extension to Hierarchical Ensembles

Consider again a random code  $\mathcal{C}$  of size  $M_1 = e^{nR_1}$ , but this time, it is drawn from a different ensemble, which is in the spirit of the ensemble of the direct part of the coding theorem for the degraded broadcast channel (see, e.g., [6, Section 15.6.2]). Specifically, let  $\mathcal{U}$  be a finite alphabet, let  $P_U$  be a given probability distribution on  $\mathcal{U}$ , and let  $P_{X|U}$  be a given matrix of conditional probabilities of  $X$  given  $U$ . We first select, independently at random,  $M = e^{nR}$   $n$ -vectors (“cloud centers”),  $\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{M-1}$ , all under the uniform distribution over the type class  $\mathcal{T}(P_U)$ . Next, for each  $w = 0, 1, \dots, M-1$ , we select conditionally independently (given  $\mathbf{u}_w$ ),  $M_2 = e^{nR_2}$  codewords,  $\mathbf{x}_{wM_2}, \mathbf{x}_{wM_2+1}, \dots, \mathbf{x}_{(w+1)M_2-1}$ , under the uniform distribution across the conditional type class  $\mathcal{T}(P_{X|U}|\mathbf{u}_w)$ . Obviously, the ensemble considered in the previous sections is a special case, where  $P_U$  is a degenerate distribution, putting all its mass on one (arbitrary) letter of  $\mathcal{U}$ . All other quantities are defined similarly as before.

We next present a more general formula of  $E^*(R_1, R_2)$ , the exact random coding error exponent of decoder (1), that accommodates the above defined ensemble. This is then the exact random coding error exponent of the optimal decoder for the weak user in the degraded broadcast channel.

Here, we no longer claim that  $E^*(R_1, R_2)$  is independent of  $R_2$  and that it is achieved by decoder (3) as well.

To present the formula of  $E^*(R_1, R_2)$ , we first need a few definitions. For a given generic joint distribution  $Q_{UXY}$ , of the random variables  $U$ ,  $X$ , and  $Y$ , let  $I_Q(X; Y|U)$  denote the conditional mutual information between  $X$  and  $Y$  given  $U$ . For a given marginal  $Q_{UY}$  of  $(U, Y)$ , let  $\mathcal{S}(Q_{UY})$  denote the set of conditional distributions  $\{Q_{X|UY}\}$  such that  $\sum_y Q_{UY}(u, y)Q_{X|UY}(x|u, y) = P_{UX}(u, x)$  for every  $(u, x) \in \mathcal{U} \times \mathcal{X}$ , where  $P_{UX} = P_U \times P_{X|U}$ . We first define

$$E_1(s, Q_{UY}) = \min_{Q_{X|UY} \in \mathcal{S}(Q_{UY})} \{[I_Q(X; Y|U) - R_2]_+ : f(Q_{XY}) + [R_2 - I_Q(X; Y|U)]_+ \geq s\}, \quad (62)$$

where  $s$  is an arbitrary real number. Next, for a given marginal  $Q_Y$ , define

$$E_2(s, Q_Y) = \min_{Q_{U|Y}} [I_Q(U; Y) + E_1(s, Q_{UY})], \quad (63)$$

where the minimization is across all  $\{Q_{U|Y}\}$  such that  $\sum_y Q_Y(y)Q_{U|Y}(u|y) = P_U(u)$  for every  $u \in \mathcal{U}$ . Finally, let

$$s_0(Q_{U_0Y}) = R_2 + \max_{\{Q_{X|U_0Y} \in \mathcal{S}(Q_{U_0Y}) : I_Q(X; Y|U_0) \leq R_2\}} [f(Q_{XY}) - I_Q(X; Y|U_0)], \quad (64)$$

and

$$s_1(Q_{U_0X_0Y}) = \max\{s_0(Q_{U_0Y}), f(Q_{X_0Y})\}. \quad (65)$$

Our extended formula for  $E^*(R_1, R_2)$  is given in the following theorem.

**Theorem 2** *Let  $R_1$  and  $R_2$  be given ( $R_2 \leq R_1$ ) and let the ensemble of codes be defined as in the first paragraph of this section. Then,*

$$E^*(R_1, R_2) = \min_{Q_{Y|X_0U_0}} \{D(Q_{Y|X_0U_0} \| P_{Y|X_0} | P_{U_0X_0}) + [E_2(s_1(Q_{U_0X_0Y}), Q_Y) - R]_+\}. \quad (66)$$

where  $(U_0, X_0)$  is a replica of  $(U, X)$ , i.e.,  $P_{U_0X_0} = P_{UX}$ .

*Proof Outline.* The proof of Theorem 2 is quite a straightforward generalization of the proof of Theorem 1, which was given in full detail in Section 4. We will therefore give here merely an outline with highlights mostly on the differences. Once again, we start from the expression,

$$P_e \stackrel{\cdot}{=} \mathbf{E} [\min\{1, M \cdot \Pr\{P(\mathbf{Y}|\mathcal{C}_1) \geq P(\mathbf{Y}|\mathcal{C}_0)\}\}], \quad (67)$$



where this time, the expectation is w.r.t. the randomness of  $\mathbf{U}_0$ ,  $\mathcal{C}_0$  and  $\mathbf{Y}$ , with the latter being the channel output in response to the input  $\mathbf{X}_0$  (which is again, the transmitted codeword, without loss of generality). Here, for a given  $\mathbf{y}$ , the pairwise error probability,  $\Pr\{P(\mathbf{y}|\mathcal{C}_1) \geq P(\mathbf{Y}|\mathcal{C}_0)\}$ , is calculated w.r.t. the randomness of  $\mathbf{U}_1$ ,  $\mathcal{C}_1 = \{\mathbf{X}_{M_2}, \mathbf{X}_{M_2+1}, \dots, \mathbf{X}_{2M_2-1}\}$ , but for a given  $\mathbf{u}_0$ , and  $\mathcal{C}_0$ .<sup>3</sup>

Defining  $s$  as in the proof of Theorem 1, the pairwise error probability is calculated once again, using the large deviations properties of  $N_1(\cdot)$ , which are now binomial random variables given  $\mathbf{u}_1$ . Thus, we first calculate the pairwise error probability conditioned on  $\mathbf{U}_1 = \mathbf{u}_1$ , and then average over  $\mathbf{U}_1$ . Now, for a given  $Q_{UXY}$ , designating the joint empirical distribution of a randomly chosen  $\mathbf{x}$  together with  $(\mathbf{u}_1, \mathbf{y})$ , the binomial random variable  $N_1(Q_{UXY})$  has  $e^{nR_2}$  trials and probability of success which is of the exponential order of  $e^{-nI_Q(X;Y|U)}$ . Everything else in this large deviations analysis remains intact. Thus,  $E_0(Q_{XY})$ , in the proof of Theorem 1, should be replaced by  $E_0(Q_{UXY})$ , which is defined by

$$E_0(Q_{UXY}) = \begin{cases} [I_Q(X;Y|U) - R_2]_+ & f(Q_{XY}) \geq s - [R_2 - I_Q(X;Y|U)]_+ \\ \infty & f(Q_{XY}) < s - [R_2 - I_Q(X;Y|U)]_+ \end{cases} \quad (68)$$

Therefore,  $E_1(s, Q_{UY})$  of the proof of Theorem 1, should now be replaced by  $E_1(s, Q_{UY})$  as defined eq. (62). The conditional pairwise error probability, that includes also conditioning on  $\mathbf{U}_1 = \mathbf{u}_1$ , is then of the exponential order of  $e^{-nE_1(s, Q_{UY})}$ . After averaging this exponential function w.r.t. the randomness of  $\mathbf{U}_1$  (thus relaxing the conditioning  $\mathbf{U}_1 = \mathbf{u}_1$ ), the resulting expression becomes of the exponential order of  $e^{-nE_2(s, Q_Y)}$ , where  $E_2(s, Q_Y)$  is defined as in (63). The remaining part of the proof is exactly in the footsteps of the proof of Theorem 1, except that here, the simplifications given near the end of the proof do not seem to hold anymore.  $\square$

## References

- [1] A. B. Carleial, "Interference channels," *IEEE Trans, Inform. Theory*, vol. IT-24, no. 1, pp. 60–70, January 1978.
- [2] T. M. Cover, "Broadcast channels," *IEEE Trans, Inform. Theory*, vol. IT-18, no. 1, pp. 2–14, January 1972.

---

<sup>3</sup>The reason that the randomness of  $\mathbf{U}_1$  is accommodated already at this stage, is that this way, the  $M - 1$  pairwise error events are all independent (given  $\mathbf{U}_0$ ,  $\mathcal{C}_0$  and  $\mathbf{y}$ ) and have identical probabilities, and so, the truncated union bound,  $\min\{1, M \cdot \Pr\{P(\mathbf{y}|\mathcal{C}_1) \geq P(\mathbf{y}|\mathcal{C}_0)\}\}$ , remains exponentially tight, as before.

- [3] T. M. Cover, “Comments on broadcast channels,” *IEEE Trans, Inform. Theory*, vol. 44. no. 6 (special commemorative issue), pp. 2524–2530, October 1998.
- [4] T. M. Cover and A. El Gamal, “Capacity theorems for the relay channel,” *IEEE Trans, Inform. Theory*, vol. IT-25. no. 5, pp. 572–584, September 1979.
- [5] T. M. Cover and C. S. K. Leung, “An achievable rate region for the multiple-access channel with feedback,” *IEEE Trans, Inform. Theory*, vol. IT-27. no. 3, pp. 292–298, May 1981.
- [6] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, John Wiley & Sons, Second Edition, New Jersey, 2006.
- [7] I. Csiszár and J. Körner, “Broadcast channels with confidential messages,” *IEEE Trans, Inform. Theory*, vol. IT-24, no. 3, pp. 339–348, May 1978.
- [8] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Academic Press, 1981.
- [9] S. I. Gel’fand and M. S. Pinsker, “Coding for channel with random parameters,” *Problems of Information and Control*, vol. 9, no. 1, pp. 19–31, 1980.
- [10] A. J. Grant, B. Rimoldi, R. L. Urbanke, and P. A. Whiting, “Rate-splitting multiple access for discrete memoryless channels,” *IEEE Trans, Inform. Theory*, vol. 47. no. 3, pp. 873–890, March 2001.
- [11] Y. Kaspi and N. Merhav, “Error exponents for broadcast channels with degraded message sets,” *IEEE Trans. Inform. Theory*, vol. 57, no. 1, pp. 101–123, January 2011.
- [12] L. H. Ozarow, “The capacity of the white Gaussian multiple access channel with feedback,” *IEEE Trans. Inform. Theory*, vol. IT-30, no. 4, pp. 623–629, July 1984.
- [13] N. Merhav, “Statistical physics and information theory,” (invited paper) *Foundations and Trends in Communications and Information Theory*, vol. 6, nos. 1–2, pp. 1–212, 2009.
- [14] N. Merhav, “Exact correct-decoding exponent for the wiretap channel decoder,” submitted to *IEEE Trans. Inform. Theory*, March 2014. <http://arxiv.org/pdf/1403.6143.pdf>
- [15] N. Shulman, *Communication over an Unknown Channel via Common Broadcasting*, Ph.D. dissertation, Department of Electrical Engineering – Systems, Tel Aviv University, July 2003. [http://www.eng.tau.ac.il/~shulman/papers/Nadav\\_PhD.pdf](http://www.eng.tau.ac.il/~shulman/papers/Nadav_PhD.pdf)

- [16] A. Somekh-Baruch and N. Merhav, “Achievable error exponents for the private fingerprinting game,” *IEEE Trans. Inform. Theory*, vol. 53, no. 5, pp. 1827–1838, May 2007.
- [17] A. D. Wyner, “The wire-tap channel,” *Bell System Technical Journal*, vol. 54, no. 8, pp. 1355–1387, Oct. 1975.
- [18] A. D. Wyner and J. Ziv, “The rate-distortion function for source coding with side information at the decoder,” *IEEE Trans. Inform. Theory*, vol. IT-22, no. 1, pp. 1–10, January 1976.