

Universal Decoding for Source–Channel Coding with Side Information ^{*}

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 a setting of Slepian–Wolf coding, where the random bin of the source vector undergoes channel coding, and then decoded at the receiver, based on additional side information, correlated to the source. For a given distribution of the randomly selected channel codewords, we propose a universal decoder that depends on the statistics of neither the correlated sources nor the channel, assuming first that they are both memoryless. Exact analysis of the random–binning/random–coding error exponent of this universal decoder shows that it is the same as the one achieved by the optimal maximum a–posteriori (MAP) decoder. Previously known results on universal Slepian–Wolf source decoding, universal channel decoding, and universal source–channel decoding, are all obtained as special cases of this result. Subsequently, we further generalize the results in several directions, including: (i) finite–state sources and finite–state channels, along with a universal decoding metric that is based on Lempel–Ziv parsing, (ii) arbitrary sources and channels, where the universal decoding is with respect to a given class of decoding metrics, and (iii) full (symmetric) Slepian–Wolf coding, where both source streams are separately fed into random–binning source encoders, followed by random channel encoders, which are then jointly decoded by a universal decoder.

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

1 Introduction

Universal decoding for unknown channels is a topic that attracted considerable attention throughout the last four decades. In [10], Goppa was the first to offer the *maximum mutual information* (MMI) decoder, which decodes the message as the one whose codeword has the largest empirical mutual information with the channel output sequence. Goppa proved that for discrete memoryless channels (DMC's), MMI decoding attains capacity. Csiszár and Körner [4, Theorem 5.2] have further showed that the random coding error exponent of the MMI decoder, pertaining to the ensemble of the uniform random coding distribution over a certain type class, achieves the same random coding error exponent as the optimum, maximum likelihood (ML) decoder. Ever since these early works on universal channel decoding, a considerably large volume of research work has been done, see, e.g., [3], [7], [8], [13], [14], [15], [18], [26], for a non-exhaustive list of works on memoryless channels, as well as more general classes of channels.

At the same time, considering the analogy between channel coding and Slepian–Wolf (SW) source coding, it is not surprising that universal schemes for SW decoding, like the *minimum entropy* (ME) decoder, have also been derived, first, by Csiszár and Körner [4, Exercise 3.1.6], and later further developed by others in various directions, see, e.g., [1], [6], [12], [27], [28], [32].

Much less attention, however, has been devoted to universal decoding for joint source–channel coding, where both the source and the channel are unknown to the decoder. Csiszár [2] was the first to propose such a universal decoder, which he referred to as the *generalized MMI* decoder. The generalized MMI decoding metric, to be maximized among all messages, is essentially¹ given by the difference between the empirical input–output mutual information of the channel and the empirical entropy of the source. In a way, it naturally combines the concepts of MMI channel decoding and ME source decoding. But the emphasis in [2] was inclined much more towards upper and lower bounds on the reliability function, whereas the universality of the decoder was quite a secondary issue. Consequently, later articles that refer to [2] also focus, first and foremost, on the joint source–channel reliability function and not really on universal decoding. We are not aware of subsequent works on universal source–channel decoding other than [17], which concerns a completely different setting, of zero-delay coding.

In this work, we consider universal joint source–channel decoding in several settings that are all more general than that of [2]. In particular, we begin by considering the communication system depicted in Fig. 1, which is described as follows: A source vector \mathbf{u} , emerging from a discrete memoryless source (DMS), undergoes Slepian–Wolf encoding (random binning) at rate R , followed by channel coding (random coding). The discrete memoryless channel (DMC) output \mathbf{y} is fed into the decoder, along with a side information (SI) vector \mathbf{v} , correlated to the source \mathbf{u} , and the output of the decoder, $\hat{\mathbf{u}}$, should agree with \mathbf{u} with probability as high as possible.

Our first step is to characterize the exact exponential rate of the expected error probability, associated with the optimum MAP decoder, where the expectation is over both ensembles of the

¹Strictly speaking, Csiszár's decoding metric is slightly different, but is asymptotically equivalent to this definition.

random binning encoder and the random channel code. We refer to this exponential rate as the *random-binning/random-coding error exponent*. The second step, which is the more important one for our purposes, is to show that this error exponent is also achieved by a universal decoder, that depends neither on the statistics of the source, nor of the channel, and which is similar to Csiszár's generalized MMI decoder. Beyond the fact this model is more general than the one in [2] (in the sense of including the random binning component as well as decoder SI), the assertion of the universal optimality of the generalized MMI decoder is stronger here than in [2]. In [2] the performance of the generalized MMI decoder is compared directly to an upper bound on the joint source-channel reliability function, and the claim on the optimality of this decoder is asserted only in the range where this bound is tight. Here, on the other hand (similarly as in earlier works on universal pure channel decoding), we argue that the generalized MMI decoder is *always* asymptotically as good as the optimal MAP decoder, in the error exponent sense, no matter whether or not there is a gap between the achievable exponent and the upper bound on the reliability function. In other words, like in previous works on universal decoding, the focus is on asymptotic optimality of the decoder for the average code and for an unknown channel, rather than on optimality of the overall communication system. However, as we shall see later on, since full optimization of the random coding ensemble is infeasible, due to channel uncertainty, the best one can hope for is the MAP source-channel error exponent due to Gallager [9, Problem 5.16]. We also provide an upper bound to the error exponent for any communication system with the configuration of Fig. 1 and discuss the conditions under which it is met.

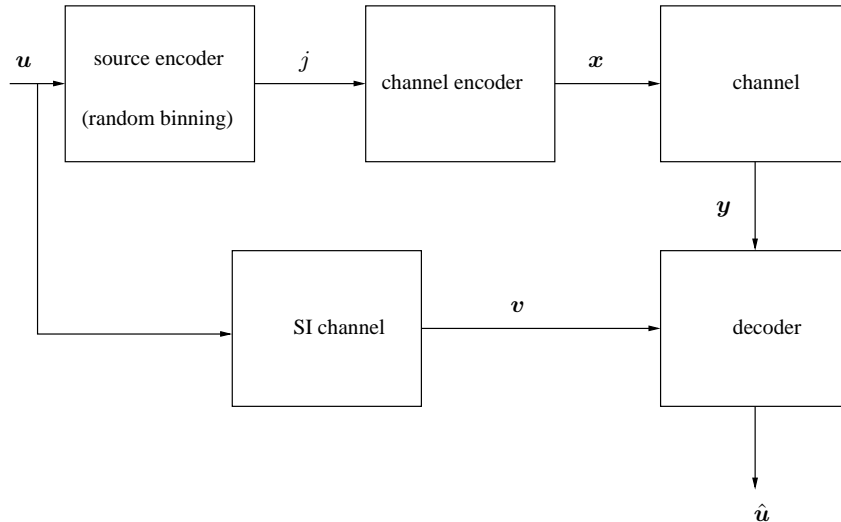


Figure 1: Slepian-Wolf source coding, followed by channel coding. The source \mathbf{u} is source-channel encoded, whereas the correlated SI \mathbf{v} (described as being generated by a DMC fed by \mathbf{u}) is available at the decoder.

One motivation for studying this model is that it captures, in a unified framework, several

important special cases of communication systems, from the perspective of universal decoding.

1. Separate source coding and channel coding without SI: letting \mathbf{v} be degenerate.
2. Pure SW source coding: letting the channel be clean ($\mathbf{y} \equiv \mathbf{x}$) and assuming the channel alphabet to be very large (so that the probability for two or more identical codewords would be negligible).
3. Pure channel coding: letting the source be binary and symmetric, and the SI be degenerate.
4. Joint source–channel coding with and without SI: letting the binning rate R be sufficiently large, so that probability of ambiguous binning (i.e., when two or more source vectors are mapped into the same bin) is negligible. In this case, the mapping between source vectors and channel input vectors is one–to–one with high probability, and therefore, this is a joint source–channel code. More details on this aspect will follow in the sequel.
5. Systematic coding: letting the SI channel (from \mathbf{u} to \mathbf{v}) in Fig. 1 be identical to the main channel (from \mathbf{x} to \mathbf{y}), and then the SI channel may represent transmission of the systematic (uncoded) part of the code (see discussions on this point of view also in [20] and [30]).

Another motivation is that it serves as the basis for the more important part of the paper, where we provide three further extensions of this communication system model. In at least two of these more general situations, the analysis is more tricky and several difficulties that are encountered need to be handled with care. The extended scenarios are the following.

1. Extending the scope from memoryless sources and channels to finite–state sources and finite–state channels. Here, the universal joint source–channel decoding metric is based on Lempel–Ziv (LZ) parsing, with the inspiration of [33]. The non–trivial parts of the analysis (not encountered in [33] or other related works) are mainly those described in items 1, 7 and 8 of Subsection 5.1.
2. Further extending the scope to arbitrary sources and channels, but allowing a given, limited class of reference decoding metrics. We propose a universal joint source–channel decoder with the property that, no matter what the underlying source and channel may be, this universal decoder is asymptotically as good as the best decoder in the class for these source and channel. This extends the recent study in [25], from pure channel coding to joint source–channel coding.
3. Generalizing to the model to separate encodings (source binning followed by channel coding) and joint decoding of two correlated sources (see Fig. 2 in Section 5.3). Here the universal decoder must handle several types of error events due to possible ambiguities in the binning encoder. As a consequence of this fact, the proposed universal decoding metric for this scenario is surprisingly different from what one may expect.

Finally, a few words are in order concerning the error exponent analysis. The ensemble of codes in our setting combines random binning (for the source coding part) and random coding (for the channel coding part), which is considerably more involved than ordinary error exponent analyses that is associated with either one but not both. This requires a rather careful analysis, in two steps,

where in the first, we take the average probability of error over the ensemble of random binning codes, for a given channel code, and at the second step, we average over the ensemble of channel codes. The latter employs the type class enumeration method [21, Chap. 6], which has already proved rather useful as a tool for obtaining exponentially tight random coding bounds in various contexts (see, e.g., [22], [23], [24] for a sample), and this work is no exception in that respect, as the resulting error exponents are tight for the average code.

The remaining part of the paper is organized as follows. In Section 2, we establish notation conventions, formalize the model and the problem, and finally, review some preliminaries. Section 3 provides the main result along with some discussion. The proof of this result appears in Section 4, and finally, Section 5 is devoted for the various extensions described above.

2 Notation Conventions, Problem Setting and Preliminaries

2.1 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. Sources and channels will be denoted by the letter P , Q , or W , subscripted by the names of the relevant random variables/vectors and their conditionings, if applicable, following the standard notation conventions, e.g., Q_X , $P_{Y|X}$, and so on. When there is no room for ambiguity, these subscripts will be omitted. To avoid cumbersome notation, the various probability distributions will be denoted as above, no matter whether probabilities of single symbols or n -vectors are addressed. Thus, for example, $P_U(u)$ (or $P(u)$) will denote the probability of a single symbol $u \in \mathcal{U}$, whereas $P_U(\mathbf{u})$ (or $P(\mathbf{u})$) will stand for the probability of the n -vector $\mathbf{u} \in \mathcal{U}^n$. The probability of an event \mathcal{E} will be denoted by $\Pr\{\mathcal{E}\}$, and the expectation operator with respect to (w.r.t.) a probability distribution P will be denoted by $\mathbf{E}\{\cdot\}$. The entropy of a generic distribution Q on \mathcal{X} will be denoted by $\mathcal{H}(Q)$. 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 is, $\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{a_n}{b_n} = 0$. Accordingly, the notation $a_n \doteq 2^{-n\infty}$ means that a_n decays at a super-exponential rate (e.g., double-exponentially). Unless specified otherwise, logarithms and exponents, throughout this paper, should be understood to be taken to the base 2. The indicator function of an event \mathcal{E} will be denoted by $\mathcal{I}\{\mathcal{E}\}$. The notation $[x]_+$ will stand for $\max\{0, x\}$. The minimum between two reals, a and b , will frequently be denoted by $a \wedge b$. The cardinality of a finite set, say \mathcal{A} , will be denoted by $|\mathcal{A}|$.

The empirical distribution of a sequence $\mathbf{x} \in \mathcal{X}^n$, which will be denoted by $\hat{P}_{\mathbf{x}}$, is the vector of

relative frequencies $\hat{P}_{\mathbf{x}}(x)$ of each symbol $x \in \mathcal{X}$ in \mathbf{x} . The type class of $\mathbf{x} \in \mathcal{X}^n$, denoted $\mathcal{T}(\mathbf{x})$, is the set of all vectors \mathbf{x}' with $\hat{P}_{\mathbf{x}'} = \hat{P}_{\mathbf{x}}$. When we wish to emphasize the dependence of the type class on the empirical distribution, say Q , we will denote it by $\mathcal{T}(Q)$. Information measures associated with empirical distributions will be denoted with ‘hats’ and will be subscripted by the sequences from which they are induced. For example, the entropy associated with $\hat{P}_{\mathbf{x}}$, which is the empirical entropy of \mathbf{x} , will be denoted² by $\hat{H}_{\mathbf{x}}(X)$. Again, the subscript will be omitted whenever it is clear from the context what sequence the empirical distribution was extracted from. Similar conventions will apply to the joint empirical distribution, the joint type class, the conditional empirical distributions and the conditional type classes associated with pairs of sequences of length n . Accordingly, $\hat{P}_{\mathbf{x}\mathbf{y}}$ would be the joint empirical distribution of $(\mathbf{x}, \mathbf{y}) = \{(x_i, y_i)\}_{i=1}^n$, $\mathcal{T}(\mathbf{x}, \mathbf{y})$ or $\mathcal{T}(\hat{P}_{\mathbf{x}\mathbf{y}})$ will denote the joint type class of (\mathbf{x}, \mathbf{y}) , $\mathcal{T}(\mathbf{x}|\mathbf{y})$ will stand for the conditional type class of \mathbf{x} given \mathbf{y} , $\hat{H}_{\mathbf{x}\mathbf{y}}(X, Y)$ will designate the empirical joint entropy of \mathbf{x} and \mathbf{y} , $\hat{H}_{\mathbf{x}\mathbf{y}}(X|Y)$ will be the empirical conditional entropy, $\hat{I}_{\mathbf{x}\mathbf{y}}(X; Y)$ will denote empirical mutual information, and so on.

2.2 Problem Setting for the Basic Setting

Let $(\mathbf{U}, \mathbf{V}) = \{(U_t, V_t)\}_{t=1}^n$ be n independent copies of a pair of random variables, $(U, V) \sim P_{UV}$, taking on values in finite alphabets, \mathcal{U} and \mathcal{V} , respectively. The vector \mathbf{U} will designate the source vector to be encoded, whereas the vector \mathbf{V} will serve as correlated SI, available to the decoder. Let W designate a DMC, with single-letter, input-output transition probabilities $W(y|x)$, $x \in \mathcal{X}$, $y \in \mathcal{Y}$, \mathcal{X} and \mathcal{Y} being finite input and output alphabets, respectively. When the channel is fed by an input vector $\mathbf{x} \in \mathcal{X}^n$, it produces³ a channel output vector $\mathbf{y} \in \mathcal{Y}^n$, according to

$$W(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^n W(y_t|x_t). \quad (1)$$

Consider the communication system depicted in Fig. 1. When a given realization $\mathbf{u} = (u_1, \dots, u_n)$, of the source vector \mathbf{U} , is fed into the system, it is encoded into one out of $M = 2^{nR}$ bins, selected independently at random for every member of \mathcal{U}^n . Here, $R > 0$ is referred to as the *binning rate*. The bin index $j = f(\mathbf{u})$ is mapped into a channel input vector $\mathbf{x}(j) \in \mathcal{X}^n$, which in turn is transmitted across the channel W . The various codewords $\{\mathbf{x}(j)\}_{j=1}^M$ are selected independently at random under the uniform distribution across a given type class $\mathcal{T}(Q)$, Q being a given probability distribution over \mathcal{X} .⁴ The randomly chosen codebook $\{\mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(M)\}$ will be denoted by \mathcal{C} . Both the channel encoder, \mathcal{C} , and the realization of the random binning source encoder, f ,

²Note that here we use the letter H in the ordinary font, as opposed to the earlier defined notation of the entropy as a functional of a distribution, where we used the calligraphic \mathcal{H} .

³Without essential loss of generality, and similarly as in [2], we assume that the source and the channel operate at the same rate, so that while the source emits the n -vector (\mathbf{U}, \mathbf{V}) , the channel is used n times exactly, transforming $\mathbf{x} \in \mathcal{X}^n$ to $\mathbf{y} \in \mathcal{Y}^n$. The extension to the case where operation rates are different (bandwidth expansion factor different from 1) is straightforward but is avoided here, in the quest of keeping notation and expressions less cumbersome.

⁴Rather than the same type $\mathcal{T}(Q)$ for all bins, a more general ensemble may allow different types of codewords to bins of different types of source vectors. We will address this point in the next section.

are revealed to the decoder as well. With a slight abuse of notation, we will sometimes denote $\mathbf{x}(j) = \mathbf{x}[f(\mathbf{u})]$ by $\mathbf{x}[\mathbf{u}]$. The optimal (MAP) decoder estimates \mathbf{u} , using the channel output $\mathbf{y} = (y_1, \dots, y_n)$ and the SI vector $\mathbf{v} = (v_1, \dots, v_n)$, according to:

$$\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} P(\mathbf{u}, \mathbf{v}) W(\mathbf{y} | \mathbf{x}[\mathbf{u}]). \quad (2)$$

The average probability of error, \overline{P}_e , is the probability of the event $\{\hat{\mathbf{U}} \neq \mathbf{U}\}$, where in addition to the randomness of (\mathbf{U}, \mathbf{V}) and the channel output \mathbf{Y} , the randomness of the source binning code and the channel code are also taken into account. The random-binning/random-coding error exponent, associated with the optimal, MAP decoder, is defined as

$$E(R, Q) = \lim_{n \rightarrow \infty} \left[-\frac{\log \overline{P}_e}{n} \right], \quad (3)$$

provided that the limit exists (a fact that will become evident from the analysis in the sequel).

The first step is to derive a single-letter expression for the exact random-binning/random-coding error exponent $E(R, Q)$. While the MAP decoder depends on the source P and the channel W , the second step is to propose a universal decoder, independent of P and W , whose average error probability decays exponentially at the same rate, $E(R, Q)$. Note that we are considering a fixed Q without attempt to maximize $E(R, Q)$ w.r.t. Q since the maximizing Q normally depends on the unknown channel W (more on this in the next subsection). Finally, our main goal is to extend the scope beyond memoryless systems, as well as to the setting where the role of \mathbf{V} is no longer merely to serve as SI at the decoder, but rather as another source vector, encoded similarly, but separately from \mathbf{U} (see Fig. 2).

2.3 Preliminaries – the Joint Source–Channel Error Exponent

To the best of our knowledge, the first to consider error exponents for joint source–channel coding (without SI) was Gallager (see also Jelinek [11]). In the second part of Problem 5.16 in his textbook [9] (pp. 534–535), the reader is requested to prove that for a given DMS P , a given DMC W , and a given product distribution Q for random selection of a channel input vector $\mathbf{x}[\mathbf{u}]$ for each source vector \mathbf{u} , the average probability of error is upper bounded by

$$\overline{P}_e \leq \exp \left\{ -n \max_{0 \leq \rho \leq 1} \left[E_0(\rho, Q) - (1 + \rho) \ln \left(\sum_{u \in \mathcal{U}} [P(u)]^{1/(1+\rho)} \right) \right] \right\}, \quad (4)$$

where $E_0(\rho, Q)$ is the well-known Gallager function

$$E_0(\rho, Q) = -\ln \sum_{y \in \mathcal{Y}} \left[\sum_{x \in \mathcal{X}} Q(x) W(y|x)^{1/(1+\rho)} \right]^{1+\rho}. \quad (5)$$

It is easy to show (see Appendix) that this exponential upper bound is equivalent to

$$\overline{P}_e \leq \exp \left\{ -n \min_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} \left[E^s(\tilde{R}) + E_r^c(\tilde{R}, Q) \right] \right\}, \quad (6)$$

where $E^s(\tilde{R})$ is the source reliability function [16], given by

$$E^s(\tilde{R}) = \min_{\{P': \mathcal{H}(P') \geq \tilde{R}\}} D(P' \| P), \quad (7)$$

$D(P' \| P)$ being the Kullback–Leibler divergence between P' and P , and

$$E_r^c(\tilde{R}, Q) = \max_{0 \leq \rho \leq 1} [E_0(\rho, Q) - \rho \tilde{R}]. \quad (8)$$

Slightly more than a decade later, Csiszár [2] derived upper and lower bounds on the reliability function of lossless joint source–channel coding (again, without SI). Csiszár has shown, in that paper, that the reliability function, E_{jsc} , of lossless joint source–channel coding is upper bounded by

$$E_{\text{jsc}} \leq \min_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} [E^s(\tilde{R}) + E^c(\tilde{R})] \quad (9)$$

where $E^c(\tilde{R})$ is the channel reliability function, for which there is a closed–form expression available only at rate zero (the zero–rate expurgated exponent) and at rates above the critical rate (the sphere–packing exponent). The lower bound in [2] is given by

$$E_{\text{jsc}} \geq \min_{\tilde{R}} [E^s(\tilde{R}) + E_r^c(\tilde{R})], \quad (10)$$

where $E_r^c(\tilde{R})$ is the random coding error exponent of the channel W and where we have relaxed the constraint on the range of \tilde{R} since this unconstrained minimum is attained in that range anyway (see [2, p. 323, one line before the Remark]). The upper and the lower bounds coincide (and hence provide the exact reliability function) whenever the minimizing \tilde{R} , of the upper bound (9), exceeds the critical rate of the channel W .

Note that the difference between the achievable exponents of Gallager and Csiszár is in the channel error exponent terms. In the former, it is $E_r^c(\tilde{R}, Q)$, whereas in the latter it is improved to $E_r^c(\tilde{R}) = \max_Q E^c(\tilde{R}, Q)$. The reason is that, while Gallager uses the same type random coding distribution Q for all codewords $\{\mathbf{x}[\mathbf{u}]\}$, Csiszár partitions the source space according to the various types $\{P'\}$ and maps each such type into a channel subcode whose rate is essentially $\tilde{R} = \mathcal{H}(P')$ and for which the channel input type Q is optimized according to $\max_Q E_r^c(\mathcal{H}(P'), Q)$. The difference disappears, of course, if for the channel W , the same Q maximizes $E^c(\tilde{R}, Q)$ for every \tilde{R} . This happens, for example, for the modulo–additive channel $Y = X \oplus N$ (N being independent of X), where the uniform distribution Q is optimal independently of \tilde{R} .

An expression equivalent to (10) is given by

$$E_{\text{jsc}} \geq \min_{P'} \max_Q \min_{W'} \{D(P' \| P) + D(W' \| W | Q) + [I(X; Y') - H(U')]\}_+ \quad (11)$$

where U' is an auxiliary random variable drawn by a source P' over \mathcal{U} (hence $H(U') = \mathcal{H}(P')$), X is governed by Q , Y' designates the output of an auxiliary channel $W' : \mathcal{X} \rightarrow \mathcal{Y}$ fed by X , and $D(W' \| W | Q)$ is the Kullback–Leibler divergence between W' and W , weighted by Q , that is

$$D(W' \| W | Q) = \sum_{x \in \mathcal{X}} Q(x) \sum_{y \in \mathcal{Y}} W'(y|x) \log \frac{W'(y|x)}{W(y|x)}. \quad (12)$$

Here, the term $D(P'\|P)$ is parallel to the source coding exponent, $E^s(\tilde{R})$, whereas the sum of other two terms can be referred to the channel coding exponent $E_r^c(\tilde{R})$ (see [2]). This is true since the minimization over P' in (11) can be carried out in two steps, where in the first, one minimizes over all $\{P'\}$ with a given entropy $\mathcal{H}(P') = \tilde{R}$ (thus giving rise to $E^s(\tilde{R})$ according to (7)), and then minimizes over \tilde{R} .

In a nutshell, the idea behind the converse part in [2] is that each type class P' , of source vectors, can be thought of as being mapped by the encoder into a separate channel subcode at rate $\tilde{R} = \mathcal{H}(P')$, and then the probability of error is lower bounded by the contribution of the worst subcode. This is to say that for the purpose of the lower bound, only decoding errors *within* each subcode are counted, whereas errors caused by confusing two source vectors that belong two different subcodes, are ignored. An interesting point, in this context, is that whenever the upper and the lower bound coincide (in the exponential scale), this means that confusions *within* the subcodes dominate the error probability, at least as far as error exponents are concerned, whereas errors of confusing codewords from different subcodes are essentially immaterial. We will witness the same phenomenon from a different perspective, in the sequel.

For the achievability part of [2], Csiszár analyzes the performance of a universal decoder, that is asymptotically equivalent to the following:

$$\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} [\hat{I}_{\mathbf{x}[\mathbf{u}]\mathbf{y}}(X; Y) - \hat{H}_{\mathbf{u}}(U)]. \quad (13)$$

As mentioned earlier, Csiszár refers to his decoder as the generalized MMI decoder. An important point to observe, however, is that in this universal setting, it makes sense to assume that the encoder does not know the channel W either and hence cannot match the optimal channel input type Q to every given source type rate $\tilde{R} = \mathcal{H}(P')$ as described above, because this optimal type depends on the unknown channel W (see [2, page 323, Remark]). Csiszár suggests, in this case, to select a fixed type Q for all source types, in which case the achievable exponent becomes the same as Gallager's exponent. Throughout this paper, we shall adopt this suggestion, and for this reason, we have defined our objective (in the previous section) to universally achieve $E(R, Q)$ for a given Q without attempt to optimize Q ,

3 Results for the System of Fig. 1

Our problem setting and results are more general than those of [2] from the following aspects: (i) we include side information \mathbf{V} , (ii) we include a cascade of random binning encoder and a channel encoder (separate source- and channel coding), and (iii) we compare the performance of the universal decoder to that of the MAP decoder (2) and show that they *always* (i.e., even when the random coding ensemble is not good enough to achieve the reliability function) have the same error exponent, whereas Csiszár compares the performance of (13) to the upper bound (9) and thus may conclude for asymptotic optimality of the decoder (together the encoder) only when the exact joint source-channel reliability function is known.

Concerning (ii), one may wonder what is the motivation for separate source– and channel coding, because joint source–channel coding is always at least as good. The answer to this question is two–fold:

1. In some applications, system constraints dictate separate source– and channel coding, for example, when the two encodings are performed at different units/locations or when general engineering considerations (like modularity) dictate separation.
2. The joint source–channel setting, for a fixed channel input type Q , can always be obtained as a special case, by choosing the binning rate R sufficiently high, since then the binning encoder is a one–to–one mapping with an overwhelmingly high probability and the channel code in cascade to the binning code is equivalent to a direct mapping between source vectors and channel input vectors.

Our main result is given by the following theorem.

Theorem 1 *Consider the problem setting defined in Subsection 2.2.*

(a) *The random–binning/random–coding error exponent of the MAP decoder is given by*

$$E(R, Q) = \min_{P_{U'V'}, W'} \{D(P_{U'V'} \| P_{UV}) + D(W' \| W|Q) + [R \wedge I(X; Y') - H(U'|V')]\}_+ \quad (14)$$

where $(U', V') \in \mathcal{U} \times \mathcal{V}$ are auxiliary random variables jointly distributed according to $P_{U'V'}$, and $Y' \in \mathcal{Y}$ is another auxiliary random variable that designates the output of channel W' when fed by $X \sim Q$.

(b) *The universal decoders,*

$$\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} [\hat{I}_{\mathbf{x}[\mathbf{u}]\mathbf{y}}(X; Y) - \hat{H}_{\mathbf{u}\mathbf{v}}(U|V)] \quad (15)$$

and

$$\hat{\mathbf{u}} = \arg \max_{\mathbf{u}} [R \wedge \hat{I}_{\mathbf{x}[\mathbf{u}]\mathbf{y}}(X; Y) - \hat{H}_{\mathbf{u}\mathbf{v}}(U|V)], \quad (16)$$

both achieve $E(R, Q)$.

Decoder (15) is, of course, a natural extension of (13) to our setting. As for (16), while it offers no apparent advantage over (15), it is given here as an alternative decoder for future reference. It will turn out later that conceptually, (16) lends itself more naturally to the extension that deals with separate encodings and joint decoding of two correlated sources, where in the extended version of (16), it will not be obvious (at least not to the author) that the operator $R \wedge (\cdot)$ is neutral (i.e., an expression like $R \wedge x$ can be simply replaced by x , as is indeed suggested here by the equivalence between (15) and (16)). Another interesting point concerning (16), is that it appears more clearly as a joint extension of the MMI decoder of pure channel decoding and the ME decoder of pure source coding. When R dominates the term $R \wedge \hat{I}_{\mathbf{x}[\mathbf{u}]\mathbf{y}}(X; Y)$, the source coding component of the problem is more prominent and (16) is essentially equivalent to the ME decoder. Otherwise, it is essentially equivalent to (15).

As can be seen, $E(R, Q)$ is monotonically non-decreasing in R ,⁵ but when R is sufficiently large, the term $R \wedge I(X; Y')$ is dominated by $I(X; Y')$ (say, $R = \log |\mathcal{X}|$), which yields saturation of $E(R, Q)$ to the level of the joint source-channel random coding exponent (for a given Q), similarly as in (11), except that here, the entropy $H(U')$ is replaced by the conditional entropy $H(U'|V')$, due to the SI. Obviously, if V' is degenerate (e.g., equal to a fixed $v \in \mathcal{V}$ with probability one), then we are back to (11). For another extreme case, if the channel is clean, Q is uniform, and the channel alphabet is very large, then $\hat{I}(X; Y') = \hat{H}(X) = \log |\mathcal{X}|$ is large as well, and then $R \wedge \hat{I}(X; Y')$ is dominated by R . In this case, we recover the SW random binning error exponent (see, e.g., [32] and references therein).

Finally, although not directly related to the aspect of universality, for the sake of completeness (and in analogy to [2]), we next provide also a converse bound that applies to any communication system of the type depicted in Fig. 1, where both the binning code $f : \mathcal{U}^n \rightarrow \{1, 2, \dots, M\}$ and the channel code, that maps $\{1, 2, \dots, M\}$ into \mathcal{C} , are arbitrary (and deterministic), and where the optimal MAP decoder is used. It is not difficult to extend Lemma 2 of [2] to the scenario under discussion and to argue that given the source P_{UV} , the channel W , and the binning rate R , the highest achievable source-channel error exponent, $E(P_{UV}, W, R)$ is upper bounded by

$$E(P_{UV}, W, R) \leq \min_{P_{U'V'}, W'} \{D(P_{U'V'} \| P_{UV}) + \mathcal{I}\{H(U'|V') \leq R\} \cdot E^c[H(U'|V')]\}, \quad (17)$$

which follows immediately from the simple consideration of viewing each conditional type $\mathcal{T}(\mathbf{u}'|\mathbf{v}')$ (whose weight is exponentially $2^{-nD(P_{U'V'} \| P_{UV})}$) as being encoded by a channel subcode at rate $\tilde{R} = H(U'|V')$, which is the corresponding empirical conditional entropy. Now, as long as $\tilde{R} \leq R$, this conditional type may be mapped into the channel code without loss of information and then the error probability within this subcode is lower bounded by $2^{-n[E^c(\tilde{R}) + o(n)]}$ (as all source messages originating from $\mathcal{T}(\mathbf{u}'|\mathbf{v}')$ are equally likely given \mathbf{v}'). However, if $\tilde{R} > R$ most of the members of $\mathcal{T}(\mathbf{u}'|\mathbf{v}')$ are mapped ambiguously by the binning encoder and the probability of error goes to unity even without the channel noise, hence the factor $\mathcal{I}\{H(U'|V') \leq R\}$ in the second term. If we further upper bound $E^c(\cdot)$ by the sphere-packing exponent, then for the second term of the above we have

$$\mathcal{I}\{H(U'|V') \leq R\} \cdot E_{\text{sp}}^c[H(U'|V')] = \max_Q \min \{D(W' \| W|Q) : R \wedge I(X; Y') \leq H(U'|V')\}. \quad (18)$$

To see why this is true, one simply examines the two cases, $H(U'|V') \leq R$ and $H(U'|V') > R$. In the former case, the constraint $R \wedge I(X; Y') \leq H(U'|V')$ is equivalent to $I(X; Y') \leq H(U'|V')$, and then both the right-hand side (r.h.s.) and the left-hand side (l.h.s.) become $E_{\text{sp}}^c[H(U'|V')]$. In the latter case, the constraint is trivially met for every W' , including the choice $W' = W$ for which the r.h.s. vanishes, exactly like the l.h.s. Putting this together, we have

$$E(P_{UV}, W, R) \leq \min_{P_{U'V'}} \max_Q \min_{\{W' : R \wedge I(X; Y') \leq H(U'|V')\}} [D(P_{U'V'} \| P_{UV}) + D(W' \| W|Q)]. \quad (19)$$

⁵ This fact is not completely trivial, since an increase in R improves the source binning part, but one may expect that it harms the channel coding part. Nonetheless, as will become apparent in the sequel (see footnote 5), the combined effect of source binning and channel coding gives a non-decreasing exponent as a function of R .

$E(R, Q)$ of Theorem 1 meets this upper bound whenever the minimizing $P_{U'V'}$ and W' of $E(R, Q)$ are such that $R \wedge I(X; Y') \leq H(U'|V')$ and that the minimization over $P_{U'V'}$ can be interchanged with the maximization over Q , e.g., when the optimal Q is independent of the coding rate, as discussed above.⁶ The first condition can be stated differently as follows: if we formally define $\tilde{E}^s(\tilde{R}) = \min\{D(P_{U'V'}\|P_{UV}) : H(U'|V') = \tilde{R}\}$ (which depends solely on the source) and $\tilde{E}^c(\tilde{R}, R, Q) = \min_{W'}\{D(W'\|W|Q) + [R \wedge I(X; Y') - \tilde{R}]_+\}$ (which depends solely on the channel), then the upper bound is attained if the value of \tilde{R} that minimizes $[E^s(\tilde{R}) + E_r^c(\tilde{R}, R, Q)]$ is large enough such that $E_r^c(\tilde{R}, R, Q)$ is achieved by W' for which $R \wedge I(X'; Y) \leq \tilde{R}$.

4 Proof of Theorem 1

The outline of the proof is as follows. We begin by showing that $E(R, Q)$ is an upper bound on the error exponent associated with the MAP decoder, and then we show that both universal decoders (15) and (16) attain $E(R, Q)$. The combination of these two facts will prove both parts of Theorem 1 at the same time.

As a first step, let the channel codebook \mathcal{C} , as well as the vectors \mathbf{u} , \mathbf{v} , $\mathbf{x} = \mathbf{x}[\mathbf{u}]$ and \mathbf{y} be given, and let $\overline{P}_e(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})$ be the average error probability given $(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})$, where the averaging is w.r.t. the ensemble of random binning source codes. For a given $\mathbf{u}' \neq \mathbf{u}$, let us define the set

$$\mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) = \mathcal{T}(Q) \cap \{\mathbf{x}' : P(\mathbf{u}', \mathbf{v})W(\mathbf{y}|\mathbf{x}') \geq P(\mathbf{u}, \mathbf{v})W(\mathbf{y}|\mathbf{x})\}. \quad (20)$$

The conditional error event, given $(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})$, is given by

$$\begin{aligned} \mathcal{E}(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C}) &= \bigcup_{\mathbf{u}' \neq \mathbf{u}} \{P(\mathbf{u}', \mathbf{v})W(\mathbf{y}|\mathbf{x}[\mathbf{u}']) \geq P(\mathbf{u}, \mathbf{v})W(\mathbf{y}|\mathbf{x}[\mathbf{u}])\} \\ &\triangleq \bigcup_{\mathbf{u}' \neq \mathbf{u}} \mathcal{E}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C}) \end{aligned} \quad (21)$$

The probability of the pairwise error event, $\mathcal{E}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})$ (again, w.r.t. the randomness of the bin assignment), is given by:

$$\overline{\Pr}\{\mathcal{E}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})\} = 2^{-nR} \left| \mathcal{C} \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right| + 2^{-nR} \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\}. \quad (22)$$

Here, the first term accounts for the probability to randomly choose a bin, other than $f(\mathbf{u})$, which is mapped to a channel input vector whose likelihood score is larger than $P(\mathbf{u}, \mathbf{v})W(\mathbf{y}|\mathbf{x}[\mathbf{u}])$. The second term is associated with the probability that $f(\mathbf{u}') = f(\mathbf{u})$ (which is 2^{-nR}), in which case the factor $W(\mathbf{y}|\mathbf{x}[\mathbf{u}']) = W(\mathbf{y}|\mathbf{x}[\mathbf{u}])$ cancels out in the pairwise likelihood score comparison, and so, \mathbf{u}' prevails if $P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})$. Now,

$$\overline{P}_e(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})$$

⁶Note that here, as opposed to [2], the independence of the optimal Q upon \tilde{R} is needed (in order to match the converse) even when the channel is *known*, because the encoder is unaware of the virtual rate $\tilde{R} = H(U'|V')$ due to the unavailability of \mathbf{v} at the encoder.

$$\begin{aligned}
&= \overline{\text{Pr}}\{\mathcal{E}(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C})\} \\
&\doteq \sum_{\{\mathcal{T}(\mathbf{u}'|\mathbf{v})\}} \overline{\text{Pr}} \left\{ \bigcup_{\mathbf{u}'' \in \mathcal{T}(\mathbf{u}'|\mathbf{v})} \mathcal{E}(\mathbf{u}, \mathbf{u}'', \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C}) \right\} \\
&\doteq \sum_{\{\mathcal{T}(\mathbf{u}'|\mathbf{v})\}} \min \left\{ 1, |\mathcal{T}(\mathbf{u}'|\mathbf{v})| \cdot 2^{-nR} \left[\left| \mathcal{C} \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right| + \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} \right] \right\}, \quad (23)
\end{aligned}$$

where we have used the fact that $\mathcal{A}(\mathbf{u}, \mathbf{u}'', \mathbf{v}, \mathbf{x}, \mathbf{y}) = \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ and $P(\mathbf{u}'', \mathbf{v}) = P(\mathbf{u}', \mathbf{v})$ for $\mathbf{u}'' \in \mathcal{T}(\mathbf{u}'|\mathbf{v})$ and where the exponential tightness of the truncated union bound (for pairwise independent events) in the last expression is known from [31, Lemma A.2, p. 109] and it can also be readily deduced from de Caen's lower bound on the probability of a union of events [5]. The next step is to average over the randomness of \mathcal{C} (except the codeword for the bin of the actual source vector \mathbf{u} , which is still given to be \mathbf{x}):

$$\begin{aligned}
\overline{P}_e(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}) &\triangleq \mathbf{E}_{\mathcal{C} \setminus \{\mathbf{x}\}} \left\{ \overline{P}_e(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}, \mathcal{C}) \right\} \\
&\doteq \sum_{\{\mathcal{T}(\mathbf{u}'|\mathbf{v})\}} \mathbf{E} \left(\min \left\{ 1, |\mathcal{T}(\mathbf{u}'|\mathbf{v})| \cdot 2^{-nR} \left[\left| \mathcal{C} \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right| + \right. \right. \right. \\
&\quad \left. \left. \left. \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} \right] \right\} \right), \quad (24)
\end{aligned}$$

where $\mathbf{E}_{\mathcal{C} \setminus \{\mathbf{x}\}}$ stands for expectation over the randomness of all codewords in \mathcal{C} other than $\mathbf{x} = \mathbf{x}[\mathbf{u}]$. Now, using the identity $\mathbf{E}\{Z\} = \int_0^\infty \Pr\{Z \geq t\} dt$, which is valid for any non-negative random variable Z , we have

$$\begin{aligned}
&\mathbf{E} \left(\min \left\{ 1, |\mathcal{T}(\mathbf{u}'|\mathbf{v})| \cdot 2^{-nR} \left[\left| \mathcal{C} \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right| + \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} \right] \right\} \right) \\
&= \int_0^1 dt \cdot \Pr \left\{ |\mathcal{T}(\mathbf{u}'|\mathbf{v})| \cdot 2^{-nR} \left[\left| \mathcal{C} \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right| + \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} \right] \geq t \right\} \\
&= \int_0^1 dt \cdot \Pr \left\{ \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} + \sum_i \mathcal{I}[\mathbf{X}(i) \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})] \geq \frac{t \cdot 2^{nR}}{|\mathcal{T}(\mathbf{u}'|\mathbf{v})|} \right\} \\
&\doteq n \ln 2 \cdot \int_0^\infty d\theta \cdot 2^{-n\theta} \Pr \left\{ \mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\} + \right. \\
&\quad \left. \sum_i \mathcal{I}[\mathbf{X}(i) \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})] \geq 2^{n[R-\theta-\hat{H}(U'|V)]} \right\}, \quad (25)
\end{aligned}$$

where in the last passage, we have used the shorthand notation $\hat{H}(U'|V)$ for $\hat{H}_{\mathbf{u}'\mathbf{v}}(U|V)$, and we have changed the integration variable from t to θ , according to the relation $t = 2^{-n\theta}$.

Consider first the case where $P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})$. Then, the integrand is given by

$$2^{-n\theta} \cdot \Pr \left\{ 1 + \sum_i \mathcal{I}[\mathbf{X}(i) \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})] \geq 2^{n[R-\theta-\hat{H}(U'|V)]} \right\} \quad (26)$$

in which the second factor is obviously equal to unity for all $\theta \geq [R - \hat{H}(U'|V)]_+$. Thus, the tail of the integral (25) is given by

$$n \ln 2 \cdot \int_{[R - \hat{H}(U'|V)]_+}^{\infty} d\theta \cdot 2^{-n\theta} \doteq 2^{-n[R - \hat{H}(U'|V)]_+}. \quad (27)$$

For $\theta < [R - \hat{H}(U'|V)]_+$, the unity term in (26) can be safely neglected, and

$$\Pr \left\{ \sum_i \mathcal{I}[\mathbf{X}(i) \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})] \geq 2^{n[R - \theta - \hat{H}(U'|V)]} \right\} \quad (28)$$

is the probability of a large deviations event associated with a binomial random variable with 2^{nR} trials and probability of success of the exponential order of 2^{-nJ} , with J being defined as

$$J \triangleq \min \left\{ \hat{I}(X'; Y) : \hat{\mathbf{P}}_{\mathbf{x}'\mathbf{y}} \text{ is such that } \mathbf{x}' \in \mathcal{T}(Q) \cap \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right\}, \quad (29)$$

where $\hat{I}(X'; Y)$ is shorthand notation for $\hat{I}_{\mathbf{x}'\mathbf{y}}(X; Y)$ and where it should be kept in mind that J depends on $\hat{\mathbf{P}}_{\mathbf{u}\mathbf{v}}$, $\hat{\mathbf{P}}_{\mathbf{u}'\mathbf{v}}$, and $\hat{\mathbf{P}}_{\mathbf{x}\mathbf{y}}$. According to [21, Chap. 6], the large deviations behavior is as follows:

$$\begin{aligned} & \Pr \left\{ \sum_i \mathcal{I}[\mathbf{X}(i) \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})] \geq 2^{n[R - \theta - \hat{H}(U'|V)]} \right\} \\ & \doteq \begin{cases} 2^{-n[J - R]_+} & R - \theta - \hat{H}(U'|V) \leq [R - J]_+ \\ 2^{-n\infty} & R - \theta - \hat{H}(U'|V) > [R - J]_+ \end{cases} \\ & = \begin{cases} 2^{-n[J - R]_+} & \theta \geq [R - \hat{H}(U'|V) - [R - J]_+]_+ \\ 2^{-n\infty} & \text{elsewhere} \end{cases} \end{aligned} \quad (30)$$

Thus, the other contribution to (25) is given by

$$\begin{aligned} & n \ln 2 \cdot \int_{[R - \hat{H}(U'|V) - [R - J]_+]_+}^{[R - \hat{H}(U'|V)]_+} d\theta \cdot 2^{-n\theta} \cdot 2^{-n[J - R]_+} \\ & \doteq \exp_2 \{ -n([R - \hat{H}(U'|V) - [R - J]_+]_+ + [J - R]_+) \} \\ & = \exp_2 \{ -n([R \wedge J - \hat{H}(U'|V)]_+ + [J - R]_+) \} \\ & = \exp_2 \{ -n([R \wedge J - \hat{H}(U'|V)]_+ - R \wedge J + R \wedge J + [J - R]_+) \} \\ & = \exp_2 \{ -n[-R \wedge J \wedge \hat{H}(U'|V) + J] \} \\ & = \exp_2 \{ -n[J - R \wedge \hat{H}(U'|V)]_+ \}, \end{aligned} \quad (31)$$

where we have repeatedly used the identity $a - [a - b]_+ = a \wedge b$. Thus, the total conditional error exponent, for the case $P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})$, is given by

$$\begin{aligned} & \min \{ [R - \hat{H}(U'|V)]_+, [J - R \wedge \hat{H}(U'|V)]_+ \} \\ & = \min \{ R - R \wedge \hat{H}(U'|V), J - R \wedge J \wedge \hat{H}(U'|V) \} \\ & = [R \wedge J - \hat{H}(U'|V)]_+, \end{aligned} \quad (32)$$

where the last line follows from the following consideration:⁷ If $\hat{H}(U'|V) > J$, then all three expressions obviously vanish and then the equality is trivially met. Otherwise, $\hat{H}(U'|V) \leq J$ implies that the term $R \wedge \hat{H}(U'|V)$, in the second line, can safely be replaced by $R \wedge J \wedge \hat{H}(U'|V)$, which makes the second line identical to $R \wedge J - R \wedge J \wedge \hat{H}(U'|V) = [R \wedge J - \hat{H}(U'|V)]_+$. In the case, $P(\mathbf{u}', \mathbf{v}) < P(\mathbf{u}, \mathbf{v})$, the conditional error exponent is just $[J - R \wedge \hat{H}(U'|V)]_+$.

Let $E_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}})$ denote the overall conditional error exponent given $(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$, i.e.,

$$E_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) = \begin{cases} [R \wedge J - \hat{H}(U'|V)]_+ & P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v}) \\ [J - R \wedge \hat{H}(U'|V)]_+ & \text{otherwise} \end{cases} \quad (33)$$

Finally, by averaging the obtained exponential estimate of $\bar{P}_e(\mathbf{U}, \mathbf{V}, \mathbf{X}, \mathbf{Y})$ over the randomness of $(\mathbf{U}, \mathbf{V}, \mathbf{X}, \mathbf{Y})$, and using the method of types in the standard manner, we obtain

$$E(R, Q) = \lim_{n \rightarrow \infty} \min_{\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}} [D(\hat{P}_{\mathbf{u}\mathbf{v}} \| P_{UV}) + D(\hat{P}_{\mathbf{y}|\mathbf{x}} \| W|Q) + E_1(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}})], \quad (34)$$

where $\hat{P}_{\mathbf{x}}$ is constrained to coincide with Q and

$$E_1(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) = \min_{\hat{P}_{\mathbf{u}'\mathbf{v}}} E_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}). \quad (35)$$

An obvious upper bound⁸ is obtained by

$$\begin{aligned} E_1(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) &\leq E_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) \\ &\leq [R \wedge \hat{I}(X; Y) - \hat{H}(U|V)]_+ \\ &\triangleq E_1^*(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}), \end{aligned} \quad (36)$$

where we have used the fact that $\mathbf{x} \in \mathcal{A}(\mathbf{u}, \mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})$ and so, for $\hat{P}_{\mathbf{u}'\mathbf{v}} = \hat{P}_{\mathbf{u}\mathbf{v}}$, one has $J \leq \hat{I}(X; Y)$. Thus,

$$\begin{aligned} E(R, Q) &\leq \min_{P_{U'V'}, W'} [D(P_{U'V'} \| P_{UV}) + D(W' \| W|Q) + E_1^*(P_{U'V'}, Q \times W')] \\ &= \min_{P_{U'V'}, W'} \left\{ D(P_{U'V'} \| P_{UV}) + D(W' \| W|Q) + [R \wedge I(X; Y') - H(U'|V')]_+ \right\} \\ &\triangleq E_U(R, Q). \end{aligned} \quad (37)$$

⁷The first line of (32) corresponds to the worst between the source coding exponent, $[R - \hat{H}(U'|V)]_+$, and the channel coding exponent, $[J - R \wedge \hat{H}(U'|V)]_+$, which is to be expected in separate source- and channel coding. While the former is non-decreasing in R , the latter is non-increasing. From the last line of (32), we learn that the overall exponent is non-decreasing in R .

⁸We are upper bounding the minimum of E_1 over $\{\hat{P}_{\mathbf{u}'\mathbf{v}}\}$ by the value of E_0 where $\hat{P}_{\mathbf{u}'\mathbf{v}} = \hat{P}_{\mathbf{u}\mathbf{v}}$, and will shortly see that this bound is actually tight. This means that the error exponent is dominated by erroneous vectors $\{\mathbf{u}'\}$ that are within the same conditional type (given \mathbf{v}) as the correct source vector \mathbf{u} . This is coherent with the observation discussed in Subsection 2.3, that errors within the subcode pertaining to the same type class dominate the error exponent.

We next argue that the universal decoder (15) achieves $E_U(R, Q)$ and hence $E(R, Q) = E_U(R, Q)$. To see why this is true, one repeats exactly the same derivation, with the following two simple modifications:

1. $\mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ is replaced by

$$\tilde{\mathcal{A}}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) = \mathcal{T}(Q) \cap \{\mathbf{x}' : \hat{I}(X'; Y) - \hat{H}(U'|V) \geq \hat{I}(X; Y) - \hat{H}(U|V)\} \quad (38)$$

and accordingly, J is replaced by

$$\tilde{J} = \min\{\hat{I}(X'; Y) : \hat{I}(X'; Y) - \hat{H}(U'|V) \geq \hat{I}(X; Y) - \hat{H}(U|V)\}. \quad (39)$$

2. The indicator function $\mathcal{I}\{P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})\}$ is replaced by $\mathcal{I}\{\hat{H}(U'|V) \leq \hat{H}(U|V)\}$.

The result is then similar except that $E_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}})$ is replaced by

$$\tilde{E}_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) = \begin{cases} [R \wedge \tilde{J} - \hat{H}(U'|V)]_+ & \hat{H}(U'|V) \leq \hat{H}(U|V) \\ [\tilde{J} - R \wedge \hat{H}(U'|V)]_+ & \text{otherwise} \end{cases} \quad (40)$$

Now, observe that for the first line of (40),

$$\begin{aligned} & R \wedge \tilde{J} - \hat{H}(U'|V) \\ & \geq R \wedge [\hat{I}(X; Y) - \hat{H}(U|V) + \hat{H}(U'|V)] - \hat{H}(U'|V) \\ & = [R - \hat{H}(U'|V)] \wedge [\hat{I}(X; Y) - \hat{H}(U|V)] \\ & \geq [R - \hat{H}(U|V)] \wedge [\hat{I}(X; Y) - \hat{H}(U|V)] \quad \text{since } \hat{H}(U'|V) \leq \hat{H}(U|V) \\ & = R \wedge \hat{I}(X; Y) - \hat{H}(U|V), \end{aligned} \quad (41)$$

where the first line follows from the definition of \tilde{J} . As for the second line,

$$\begin{aligned} \tilde{J} - R \wedge \hat{H}(U'|V) & \geq \hat{I}(X; Y) - \hat{H}(U|V) + \hat{H}(U'|V) - R \wedge \hat{H}(U'|V) \\ & = \hat{I}(X; Y) - \hat{H}(U|V) + [\hat{H}(U'|V) - R]_+ \\ & \geq \hat{I}(X; Y) - \hat{H}(U|V) + [\hat{H}(U|V) - R]_+ \quad \text{since } \hat{H}(U'|V) > \hat{H}(U|V) \\ & = \hat{I}(X; Y) - R \wedge \hat{H}(U|V) \\ & \geq R \wedge \hat{I}(X; Y) - \hat{H}(U|V). \end{aligned} \quad (42)$$

We conclude then that, no matter whether $\hat{H}(U'|V) \leq \hat{H}(U|V)$ or $\hat{H}(U'|V) > \hat{H}(U|V)$, we always have:

$$\tilde{E}_0(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{u}'\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}) \geq [R \wedge \hat{I}(X; Y) - \hat{H}(U|V)]_+ = E_1^*(\hat{P}_{\mathbf{u}\mathbf{v}}, \hat{P}_{\mathbf{x}\mathbf{y}}), \quad (43)$$

and so, the overall exponent $E_U(R)$ is achieved by (15).

As for the alternative universal decoding metric (16), the derivation is, once again, the very same, with the pairwise error event $\mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ and the variable \tilde{J} redefined accordingly as the

minimum of $\hat{I}(X'; Y)$ s.t. $R \wedge \hat{I}(X'; Y) - \hat{H}(U'|V) \geq R \wedge \hat{I}(X; Y) - \hat{H}(U|V)$, and then the last two inequalities are modified as follows: Instead of (41), we have

$$\begin{aligned}
& R \wedge \tilde{J} - \hat{H}(U'|V) \\
&= R \wedge R \wedge \hat{I}(X'; Y) - \hat{H}(U'|V) \\
&\geq R \wedge [R \wedge \hat{I}(X; Y) - \hat{H}(U|V) + \hat{H}(U'|V)] - \hat{H}(U'|V) \\
&= [R - \hat{H}(U'|V)] \wedge \{[R \wedge \hat{I}(X; Y)] - \hat{H}(U|V)\} \\
&\geq [R - \hat{H}(U|V)] \wedge \{[R \wedge \hat{I}(X; Y)] - \hat{H}(U|V)\} \quad \text{since } \hat{H}(U'|V) \leq \hat{H}(U|V) \\
&= R \wedge \hat{I}(X; Y) - \hat{H}(U|V).
\end{aligned} \tag{44}$$

and instead of (42):

$$\begin{aligned}
\tilde{J} - R \wedge \hat{H}(U'|V) &\geq R \wedge \tilde{J} - R \wedge \hat{H}(U'|V) \\
&\geq R \wedge \hat{I}(X; Y) - \hat{H}(U|V) + \hat{H}(U'|V) - R \wedge \hat{H}(U'|V) \\
&= R \wedge \hat{I}(X; Y) - \hat{H}(U|V) + [\hat{H}(U'|V) - R]_+ \\
&\geq R \wedge \hat{I}(X; Y) - \hat{H}(U|V) + [\hat{H}(U|V) - R]_+ \quad \text{since } \hat{H}(U'|V) > \hat{H}(U|V) \\
&= R \wedge \hat{I}(X; Y) - R \wedge \hat{H}(U|V) \\
&\geq R \wedge \hat{I}(X; Y) - \hat{H}(U|V).
\end{aligned} \tag{45}$$

This completes the proof of Theorem 1.

5 Extensions

As mentioned in the Introduction, in this section, we provide extensions of the above results in several directions, including: (i) finite-state sources and channels with LZ universal decoding metrics, (ii) arbitrary sources and channels with universal decoding w.r.t. a given class of metric decoders, and (iii) separate source-channel encodings and joint universal decoding of correlated source. While in (i) and (ii) we no longer expect to have single-letter formulae for the error exponent, we will still be able to propose asymptotically optimum universal decoding metrics in the error exponent sense. While the skeleton of the analysis builds upon the one of the proof of Theorem 1, we will highlight the non-trivial differences and the modifications needed relative to the proof of Theorem 1.

5.1 Finite-State Sources/Channels and a Universal LZ Decoding Metric

In [33], Ziv considered the class of finite-state channels and proposed a universal decoding metric that is based on conditional LZ parsing. Here, we discuss a similar model with a suitable extension of Ziv's decoding metric in the spirit of the generalized MMI decoder.

Consider a sequence of pairs of random variables $\{(U_i, V_i)\}_{i=1}^n$, drawn from a finite-alphabet, finite-state source, defined according to

$$P(\mathbf{u}, \mathbf{v}) = \prod_{t=1}^n P(u_t, v_t | s_t) \quad (46)$$

where s_t is the joint state of the two sources at time t , which evolves according to

$$s_t = g(s_{t-1}, u_{t-1}, v_{t-1}), \quad (47)$$

with $g : \mathcal{S} \times \mathcal{U} \times \mathcal{V} \rightarrow \mathcal{S}$ being the source next-state function, and \mathcal{S} being a finite set of states. The initial state, s_1 , is assumed to be an arbitrary fixed member of \mathcal{S} . By the same token, the channel is also assumed to be finite-state (as in [33]), i.e.,

$$W(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^n W(y_t | x_t, z_t), \quad z_t = h(z_{t-1}, x_{t-1}, y_{t-1}), \quad (48)$$

where z_t is the channel state at time t , taking on values in a finite set \mathcal{Z} and $h : \mathcal{Z} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$ is the channel next-state function. Once again, the initial state, z_1 , is an arbitrary member of \mathcal{Z} .

The remaining details of the communication system are the same as described in Subsection 2.2, with the exception that the random coding distribution, now denoted by $Q(\mathbf{x})$, is allowed here to be more general than a uniform distribution across a type class (or the uniform distribution across \mathcal{X}^n , as assumed in [33]). Similarly as in [25], we assume that Q may be any exchangeable probability distribution (i.e., \mathbf{x}' is a permutation of \mathbf{x} implies $Q(\mathbf{x}') = Q(\mathbf{x})$), and that, moreover, if the state variable z_t includes a component, say, σ_t , that is fed merely by $\{x_t\}$ (but not $\{y_t\}$), then it is enough that Q would be invariant within conditional types of \mathbf{x} given $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_n)$.

Let $\hat{H}_{\text{LZ}}(\mathbf{x}|\mathbf{y})$ denote the normalized conditional LZ compressibility of \mathbf{x} given \mathbf{y} , as defined in [33, eq. (20)] (and denoted by $u(\mathbf{x}, \mathbf{y})$ therein).⁹ Next define

$$\hat{I}_{\text{LZ}}(\mathbf{x}; \mathbf{y}) = -\frac{\log Q(\mathbf{x})}{n} - \hat{H}_{\text{LZ}}(\mathbf{x}|\mathbf{y}), \quad (49)$$

and finally, define the universal decoder

$$\tilde{\mathbf{u}} = \arg \max_{\mathbf{u}} \left[\hat{I}_{\text{LZ}}(\mathbf{x}[\mathbf{u}]; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) \right]. \quad (50)$$

Note that the first term on the r.h.s. of (49) plays a role analogous to that of the unconditional empirical entropy, $\hat{H}_{\mathbf{x}}(X)$, of the memoryless case (and indeed, at least for the uniform distribution over a type class, as assumed in the previous sections, it is asymptotically equivalent), and so, the difference in (49) makes sense as an extension of the empirical mutual information between \mathbf{x} and \mathbf{y} .

⁹This means that $n\hat{H}_{\text{LZ}}(\mathbf{x}|\mathbf{y})$ is the length of the conditional Lempel–Ziv code for \mathbf{x} , where \mathbf{y} serves as SI available to both the encoder and decoder, which is based on joint incremental parsing of the sequence pair (\mathbf{x}, \mathbf{y}) (see also [19]). Here, we are deliberately using a somewhat different notation than the usual, which hopefully makes the analogy to the memoryless case self-evident.

As in Theorem 1, part b (and as an extension to [33]), we now argue that the universal decoder (50) achieves an average error probability that is, within a sub-exponential function of n , the same as the average error probability of the MAP decoder for the given source (46) and channel (48).

Theorem 2 *Consider the problem setting defined in Subsection 2.2, with a finite-state source (46) and a finite-state channel (48). Assume that the random binning ensemble is as before and that the random channel coding distribution Q is as described in the third paragraph of this subsection. Let $\overline{P}_e^{MAP}(n)$ denote the average error probability of the MAP decoder and let $\overline{P}_e^u(n)$ denote the average error probability of the decoder (50). Then,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{\overline{P}_e^u(n)}{\overline{P}_e^{MAP}(n)} = 0. \quad (51)$$

In other words, similarly as in [33], while we do not have a characterization of the error exponent, we can still guarantee that whenever the MAP decoder has an exponentially decaying average error probability, then so does the decoder (50), and with the same exponential rate.

Proof outline. The skeleton of the proof of Theorem 2 is similar to the proof of Theorem 1, but as mentioned earlier, some non-trivial modifications are needed. Below we outline the main steps, highlighting the main modifications required.

1. The conditional type class of \mathbf{u} given \mathbf{v} , $\mathcal{T}(\mathbf{u}|\mathbf{v})$, is redefined as the set $\{\mathbf{u}' : P(\mathbf{u}', \mathbf{v}) = P(\mathbf{u}, \mathbf{v})\}$, where P is given as in (46).¹⁰ Obviously, for every given \mathbf{v} , the various ‘types’ $\{\mathcal{T}(\mathbf{u}|\mathbf{v})\}$ are equivalence classes, and hence form a partition of \mathcal{U}^n . One important property that would be essential for the proof is that the number $K_n(\mathbf{v})$ of distinct types, $\{\mathcal{T}(\mathbf{u}|\mathbf{v})\}$, under this definition, for a given \mathbf{v} , grows sub-exponentially in n (just like in the case of ordinary types). This guarantees that the probability of a union of events, over $\{\mathcal{T}(\mathbf{u}'|\mathbf{v})\}$, is of the same exponential order as the maximum term, as was the case in the proof of Theorem 1. Interestingly, this can easily be proved using the theory of LZ data compression:

$$K_n(\mathbf{v}) = \sum_{\mathbf{u} \in \mathcal{U}^n} \frac{1}{|\mathcal{T}(\mathbf{u}|\mathbf{v})|} \leq \sum_{\mathbf{u} \in \mathcal{U}^n} 2^{-n\hat{H}_{LZ}(\mathbf{u}|\mathbf{v}) + o(n)} \leq 2^{o(n)}, \quad (52)$$

where $o(n)$ stands for a sub-linear term (i.e., $\lim_{n \rightarrow \infty} o(n)/n = 0$, uniformly in both \mathbf{u} and \mathbf{v}), the first inequality is by [33, Lemma 1, p. 459]¹¹ and the second inequality is due to the fact that $n\hat{H}_{LZ}(\mathbf{u}|\mathbf{v})$ is (within negligible terms) a legitimate length function for lossless compression of \mathbf{u} (with SI \mathbf{v}) (see [33, Lemma 2] and [19]) and hence must satisfy the Kraft inequality for every given \mathbf{v} .

2. The quantity $\hat{H}(U'|V)$, in the proof of Theorem 1, is replaced by $\frac{1}{n} \log |\mathcal{T}(\mathbf{u}'|\mathbf{v})|$ with the above modified definition of the conditional type.

¹⁰Note that the requirement $P(\mathbf{u}', \mathbf{v}) = P(\mathbf{u}, \mathbf{v})$ is imposed here only for the given P , not even for every finite-state source in the class.

¹¹Not to be confused with the lemma on page 456 of [33], which is also referred to as Lemma 1.

3. The definition of J is changed to

$$J = \min \left\{ -\frac{1}{n} \log Q[\mathcal{T}(\mathbf{x}'|\mathbf{y})] : \mathbf{x}' \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right\}, \quad (53)$$

where $\mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ is the pairwise error event pertaining to the MAP decoder (for the lower bound) or to the universal decoder (50) (for the upper bound). By our assumptions, Q assigns the same probability to all members of $\mathcal{T}(\mathbf{x}'|\mathbf{y})$, thus

$$\frac{1}{n} \log Q[\mathcal{T}(\mathbf{x}'|\mathbf{y})] = \frac{\log Q(\mathbf{x}')}{n} + \frac{\log |\mathcal{T}(\mathbf{x}'|\mathbf{y})|}{n}. \quad (54)$$

4. Using the above, and following the same steps as in the proof of Theorem 1, the conditional average error probability, given $(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})$, associated with the MAP decoder, can be shown to be lower bounded by an expression of the exponential order of $\exp\{-nE_0(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})\}$, where

$$\begin{aligned} E_0(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}) &\leq \left[\min \left\{ R, -\frac{1}{n} \log Q[\mathcal{T}(\mathbf{x}|\mathbf{y})] \right\} - \frac{1}{n} \log |\mathcal{T}(\mathbf{u}|\mathbf{v})| \right]_+ \\ &= \left[\min \left\{ R, -\frac{1}{n} \log Q(\mathbf{x}) - \frac{1}{n} \log |\mathcal{T}(\mathbf{x}|\mathbf{y})| \right\} - \frac{1}{n} \log |\mathcal{T}(\mathbf{u}|\mathbf{v})| \right]_+ \\ &\leq \left[\min \left\{ R, -\frac{1}{n} \log Q(\mathbf{x}) - \hat{H}_{\text{LZ}}(\mathbf{x}|\mathbf{y}) \right\} - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) \right]_+ + o(1) \\ &= \left[R \wedge \hat{I}_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) \right]_+ + o(1) \\ &\triangleq E_1^*(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}), \end{aligned} \quad (55)$$

and where we have used twice Lemma 1 of [33, p. 459] and the fact that $\mathbf{x} \in \mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ and so, for $P(\mathbf{u}', \mathbf{v}) = P(\mathbf{u}, \mathbf{v})$, one has $J \leq -\frac{1}{n} \log Q[\mathcal{T}(\mathbf{x}|\mathbf{y})]$.

5. For the upper bound on the error probability of (50), $\mathcal{A}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})$ is replaced by

$$\tilde{\mathcal{A}}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) = \{\mathbf{x}' : I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v})\} \quad (56)$$

and accordingly, J is replaced by

$$\begin{aligned} \tilde{J} &= \min\{I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) : I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v})\} \\ &= I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v}). \end{aligned} \quad (57)$$

6. The indicator function $\mathcal{I}[P(\mathbf{u}', \mathbf{v}) \geq P(\mathbf{u}, \mathbf{v})]$ is replaced by $\mathcal{I}[\hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v}) \leq \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v})]$.

7. For the error probability analysis of the universal decoder (50), the union over erroneous source vectors $\{\mathbf{u}'\}$ is partitioned into (a sub-exponential number of) ‘types’ of the form

$$\mathcal{T}_\ell(\mathbf{u}'|\mathbf{v}) = \{\tilde{\mathbf{u}} : P(\tilde{\mathbf{u}}, \mathbf{v}) = P(\mathbf{u}', \mathbf{v}), n\hat{H}_{\text{LZ}}(\tilde{\mathbf{u}}|\mathbf{v}) = \ell\}, \quad (58)$$

for $\ell = 1, 2, \dots$, and one uses the fact that $|\mathcal{T}_\ell(\mathbf{u}'|\mathbf{v})| \leq 2^\ell$, as $n\hat{H}_{\text{LZ}}(\cdot|\mathbf{v})$ is a length function of a lossless data compression algorithm.

8. It is observed that $\sum_i \mathcal{I}[\mathbf{X}(i) \in \tilde{\mathcal{A}}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y})]$ is a binomial random variables with 2^{nR} trials and probability of success of the exponential order of $2^{-n\bar{J}}$. To see why the latter is true, consider the following:

$$\begin{aligned}
& Q \left\{ \mathbf{X}' \in \tilde{\mathcal{A}}(\mathbf{u}, \mathbf{u}', \mathbf{v}, \mathbf{x}, \mathbf{y}) \right\} \\
&= Q \left\{ I_{\text{LZ}}(\mathbf{X}'; \mathbf{y}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v}) \right\} \\
&= \sum_{\{\mathbf{x}': I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})\}} Q(\mathbf{x}') \\
&= \sum_{\{\mathbf{x}': I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})\}} Q(\mathbf{x}') 2^{n\hat{H}_{\text{LZ}}(\mathbf{x}'|\mathbf{y})} \cdot 2^{-n\hat{H}_{\text{LZ}}(\mathbf{x}'|\mathbf{y})} \\
&= \sum_{\{\mathbf{x}': I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) \geq I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})\}} \exp_2 \{ -n I_{\text{LZ}}(\mathbf{x}'; \mathbf{y}) \} \cdot 2^{-n\hat{H}_{\text{LZ}}(\mathbf{x}'|\mathbf{y})} \\
&\leq \sum_{\mathbf{x}' \in \mathcal{X}^n} \exp_2 \{ -n [I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})] \} \cdot 2^{-n\hat{H}_{\text{LZ}}(\mathbf{x}'|\mathbf{y})} \\
&\leq \exp_2 \{ -n [I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})] \} \sum_{\mathbf{x}' \in \mathcal{X}^n} 2^{-n\hat{H}_{\text{LZ}}(\mathbf{x}'|\mathbf{y})} \\
&\leq \exp_2 \{ -n [I_{\text{LZ}}(\mathbf{x}; \mathbf{y}) - \hat{H}_{\text{LZ}}(\mathbf{u}|\mathbf{v}) + \hat{H}_{\text{LZ}}(\mathbf{u}'|\mathbf{v})] + o(n) \}, \tag{59}
\end{aligned}$$

where in the last step, we have used again Kraft's inequality.

9. Using exactly the same method as in the proof of Theorem 1, one can show that that conditional error probability of the universal decoder (50) is upper bounded by an expression whose exponential order is lower bounded by $E_1^*(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})$.

It should be noted that these results continue to apply for *arbitrary* sources and channels (even deterministic ones), where the assertion would be that the decoder (50) competes favorably (in the error exponent sense) relative to any decoding metric of the form

$$\sum_{t=1}^n m_s(u_t, v_t, s_t) + \sum_{t=1}^n m_c(x_t, y_t, z_t), \tag{60}$$

where s_t and z_t evolve according to next-state functions h and g , as defined above. This follows from the observation that the assumption on underlying finite-state sources and finite-state channels was actually used merely in the assumed structure of the MAP decoding metric, with which decoder (50) competes. The fact that the overall probability of error is eventually averaged over all source vectors and channel noise realizations pertaining to finite-state probability distributions, was not really used here, since we compared the conditional error probabilities given $(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})$. The same observation has been exploited also in [25] for universal pure channel coding, and it will be further developed in the next subsection.

5.2 Arbitrary Sources and Channels With a Given Class of Metric Decoders

In [25], the following setting of universal channel decoding was studied: Given a random coding distribution Q , for independent random selection of 2^{nR} codewords $\{\mathbf{x}_i\}$, and given a limited class of reference decoders, defined by a family of decoding metrics $\{m_\theta(\mathbf{x}, \mathbf{y}), \theta \in \Theta\}$ (θ being an index or a parameter), find a decoding metric that is universal in the sense of achieving an average error probability that is, within a sub-exponential function of n , as good as the best decoder in the class, no matter what the underlying channel, $W(\mathbf{y}|\mathbf{x})$, may be. The following decoder was shown in [25] to possess this property under a certain condition that will be specified shortly: estimate the message i as the one that minimizes $u(\mathbf{x}_i, \mathbf{y}) = \log Q[\mathcal{T}(\mathbf{x}_i|\mathbf{y})]$, where $\mathcal{T}(\mathbf{x}|\mathbf{y})$ designates a notion of a “type” induced by the family of decoding metrics (rather than by channels), namely,

$$\mathcal{T}(\mathbf{x}|\mathbf{y}) = \{\mathbf{x}' : m_\theta(\mathbf{x}', \mathbf{y}) = m_\theta(\mathbf{x}, \mathbf{y}) \forall \theta \in \Theta\}. \quad (61)$$

As $\{\mathcal{T}(\mathbf{x}|\mathbf{y})\}$ are equivalence classes, they form a partition of \mathcal{X}^n for every given \mathbf{y} . The condition required for the universality of this decoding metric is that the number of distinct ‘types’ $\{\mathcal{T}(\mathbf{x}|\mathbf{y})\}$ would grow sub-exponentially with n .

A similar approach can be taken in the present problem setting. Given a family of decoding metrics of the form¹²

$$m_\theta(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}) = m_{s,\theta}(\mathbf{u}, \mathbf{v}) + m_{c,\theta}(\mathbf{x}, \mathbf{y}), \quad \theta \in \Theta, \quad (62)$$

let us define

$$\mathcal{T}_s(\mathbf{u}|\mathbf{v}) = \{\mathbf{u}' : m_{s,\theta}(\mathbf{u}', \mathbf{v}) = m_{s,\theta}(\mathbf{u}, \mathbf{v}) \forall \theta \in \Theta\} \quad (63)$$

$$\mathcal{T}_c(\mathbf{x}|\mathbf{y}) = \{\mathbf{x}' : m_{c,\theta}(\mathbf{x}', \mathbf{y}) = m_{c,\theta}(\mathbf{x}, \mathbf{y}) \forall \theta \in \Theta\}, \quad (64)$$

and assume, as before, that the numbers of distinct ‘types’, $\{\mathcal{T}_s(\mathbf{u}|\mathbf{v})\}$ and $\{\mathcal{T}_c(\mathbf{x}|\mathbf{y})\}$, both grow sub-exponentially with n . Then, the universal decoder

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u}} \{\log |\mathcal{T}_s(\mathbf{u}|\mathbf{v})| + \log Q[\mathcal{T}_c(\mathbf{x}[\mathbf{u}]|\mathbf{y})]\} \quad (65)$$

competes favorably with all metrics in the above family, no matter what the underlying source and the underlying channel may be. The proof combines the ideas of the proof of Theorem 1 above with those of [25], with the proper adjustments, of course, but it is otherwise straightforward. Here, the term $\log |\mathcal{T}_s(\mathbf{u}|\mathbf{v})|$ is the analogue of n times the conditional empirical entropy pertaining to the source part, whereas the term $\log Q[\mathcal{T}_c(\mathbf{x}[\mathbf{u}]|\mathbf{y})]$ plays the role of n times the negative empirical mutual information between $\mathbf{x}[\mathbf{u}]$ and \mathbf{y} . Therefore if, for example,

$$m_{c,\theta}(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^n m_{c,\theta}(x_t, y_t) \quad \text{and} \quad (66)$$

¹²This additive structure can be justified by the fact that the MAP decoding metric is also additive, as it maximizes $\log P(\mathbf{u}, \mathbf{v}) + \log W(\mathbf{y}|\mathbf{x}[\mathbf{u}])$.

$$m_{s,\theta}(\mathbf{u}, \mathbf{v}) = \sum_{t=1}^n m_{s,\theta}(u_t, v_t), \quad (67)$$

as is the case when the sources and the channel are memoryless, then $\{\mathcal{T}_c(\mathbf{x}|\mathbf{y})\}$ and $\{\mathcal{T}_s(\mathbf{u}|\mathbf{v})\}$ become conditional type classes in the usual sense, and we are back to the generalized MMI decoder of Section 3, provided that Q is, again, the uniform distribution within a single type class. As a final note, in this context, we mention that in this setting, the input and the output alphabets of the channel may also be continuous, see, e.g., [25, p. 5575, Example 3].

5.3 Separate Encodings and Universal Joint Decoding of Correlated Sources

Consider the system depicted in Fig. 2, which illustrates a scenario of separate source–channel encodings and joint decoding of two correlated sources, \mathbf{u}_1 and \mathbf{u}_2 . For the sake of simplicity of the presentation, we return to the assumption of memoryless systems, as in Section 3.

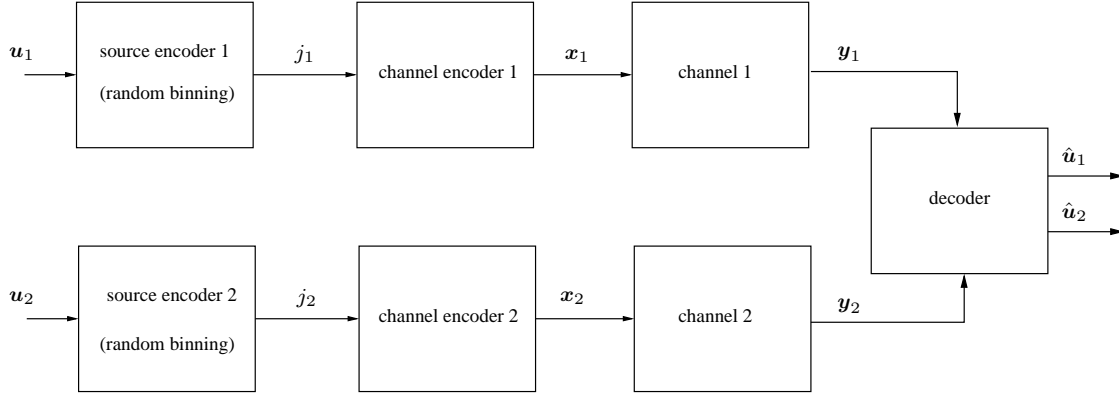


Figure 2: Separate source–channel encodings and joint decoding of two correlated sources.

Consider n independent copies $\{(U_{1,i}, U_{2,i})\}_{i=1}^n$ of a finite-alphabet pair of random variables $(U_1, U_2) \sim P_{U_1 U_2}$, as well as n uses of two independent finite-alphabet DMC's $W_1(\mathbf{y}_1|\mathbf{x}_1) = \prod_{t=1}^n W_1(y_{1,t}|x_{1,t})$ and $W_2(\mathbf{y}_2|\mathbf{x}_2) = \prod_{t=1}^n W_2(y_{2,t}|x_{2,t})$. For $k = 1, 2$, consider the following mechanism: The source vector $\mathbf{u}_k = (u_{k,1}, \dots, u_{k,n})$ is encoded into one out of $M_k = 2^{nR_k}$ bins, selected independently at random for every member of \mathcal{U}_k^n . The bin index $j_k = f_k(\mathbf{u}_k)$ is in turn mapped into a channel input vector $\mathbf{x}_k(i) \in \mathcal{X}_1^n$, which is transmitted across the channel W_k . The various codewords $\{\mathbf{x}_k(i)\}_{i=1}^{M_k}$ are selected independently at random under the uniform distribution within given type classes $\mathcal{T}(Q_k)$, where Q_k is a given distribution across \mathcal{X}_k . The randomly chosen codebook $\{\mathbf{x}_k(1), \mathbf{x}_k(2), \dots, \mathbf{x}_k(M_k)\}$ will be denoted by \mathcal{C}_k . Similarly, as before, we will sometimes denote $\mathbf{x}_k(j_k) = \mathbf{x}_k[f_k(\mathbf{u}_k)]$ by $\mathbf{x}_k[\mathbf{u}_k]$. The optimal (MAP) decoder estimates $(\mathbf{u}_1, \mathbf{u}_2)$, using the channel outputs \mathbf{y}_1 and \mathbf{y}_2 , according to

$$(\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2) = \arg \max_{\mathbf{u}_1, \mathbf{u}_2} P(\mathbf{u}_1, \mathbf{u}_2) W_1(\mathbf{y}_1|\mathbf{x}_1[\mathbf{u}_1]) W_2(\mathbf{y}_2|\mathbf{x}_2[\mathbf{u}_2]). \quad (68)$$

The main structure of the analysis continues to be essentially the same as in Section 4. The situation here, however, is significantly more involved, because five different types of pairwise error events $\{(\mathbf{u}_1, \mathbf{u}_2) \rightarrow (\mathbf{u}'_1, \mathbf{u}'_2)\}$ should be carefully handled:

1. $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 = \mathbf{u}_2$.
2. $\mathbf{u}'_2 \neq \mathbf{u}_2$ and $\mathbf{u}'_1 = \mathbf{u}_1$.
3. Both $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 \neq \mathbf{u}_2$, but (at least) \mathbf{u}'_2 is mapped into the same bin as \mathbf{u}_2 .
4. Both $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 \neq \mathbf{u}_2$, but (at least)¹³ \mathbf{u}'_1 is mapped into the same bin as \mathbf{u}_1 .
5. Both $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 \neq \mathbf{u}_2$, and neither \mathbf{u}'_1 nor \mathbf{u}'_2 belongs to the same bin as the respective true source vector.

Errors of types 1 and 2 are of the same nature as in Section 3, where the source that is estimated correctly, is actually in the role of SI at the decoder. Following (16), the respective metrics are¹⁴

$$f_1(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = R_1 \wedge \hat{I}(X_1; Y_1) - \hat{H}(U_1|U_2) \quad (69)$$

$$f_2(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = R_2 \wedge \hat{I}(X_2; Y_2) - \hat{H}(U_2|U_1). \quad (70)$$

where $\hat{I}(X_1; Y_1)$ and $\hat{H}(U_1|U_2)$ are shorthand notations for $\hat{I}_{\mathbf{x}_1 \mathbf{x}_2}(X_1; Y_1)$ and $\hat{H}_{\mathbf{u}_1 \mathbf{u}_2}(U_1|U_2)$, respectively, and so on. Errors of types 3 and 4 will turn out to be addressed by metrics of the form

$$f_3(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = R_1 \wedge \hat{I}(X_1; Y_1) + R_2 - \hat{H}(U_1, U_2) \quad (71)$$

$$f_4(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = R_2 \wedge \hat{I}(X_2; Y_2) + R_1 - \hat{H}(U_1, U_2). \quad (72)$$

Finally, error of type 5 is accommodated by

$$\begin{aligned} f_5(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) &= \hat{I}(X_1; Y_1) + \hat{I}(X_2; Y_2) - \\ &\quad \min\{\hat{H}(U_1, U_2), R_1 \wedge \hat{I}(X_1; Y_1) + R_2 \wedge \hat{I}(X_2; Y_2)\} \\ &\equiv [R_1 \wedge \hat{I}(X_1; Y_1) + R_2 \wedge \hat{I}(X_2; Y_2) - \hat{H}(U_1, U_2)]_+ + \\ &\quad [\hat{I}(X_1; Y_1) - R_1]_+ + [\hat{I}(X_2; Y_2) - R_2]_+ \end{aligned} \quad (73)$$

But we need a *single* universal decoding metric that copes with all five types of errors at the same time.

Similarly as in [25, eqs. (57)-(60)], this objective is accomplished by a metric which is given by the minimum among all five metrics above, i.e., we define our decoding metric as

$$f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = \min_{1 \leq i \leq 5} f_i(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2), \quad (74)$$

Our main result in this subsection is the following.

¹³Here, we are counting twice the case “ $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 \neq \mathbf{u}_2$ and both estimates are in the bins of their respective true source vectors.” This is done simply for symmetry the structure above, without affecting the error exponent.

¹⁴Note that f_1 does not really depend on $(\mathbf{x}_2, \mathbf{y}_2)$, and similarly, f_2 does not depend on $(\mathbf{x}_1, \mathbf{y}_1)$. Nonetheless, we deliberately adopt this uniform notation for convenience later on.

Theorem 3 Consider the above described setting of separate encodings and joint decoding of two correlated memoryless sources transmitted over two respective, independent memoryless channels. Then, the universal decoder

$$(\tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2) = \arg \max_{\mathbf{u}_1, \mathbf{u}_2} f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1[\mathbf{u}_1], \mathbf{x}_2[\mathbf{u}_2], \mathbf{y}_1, \mathbf{y}_2) \quad (75)$$

achieves the same random-binning/random-coding error exponent as the MAP decoder (68).

Proof outline. The conditional probability of error given $(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)$, for both the MAP decoder and the universal decoder, can be shown to be of the exponential order of

$$\exp_2\{-n[f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)]_+\}.$$

To show this, the analysis of the probability of error, for both the MAP decoder and the universal decoder, should be divided into several parts, according to the various types of error events. Errors of types 1 and 2 are addressed exactly as in Section 4. The more complicated part of the analysis is due to errors of types 3–5, where both competing source vectors are in error. However, this analysis too follows the same basic ideas. Here we will outline only the main ingredients that are different from those of the proof of Theorem 1.

For a given $\mathbf{u}'_1 \neq \mathbf{u}_1$ and $\mathbf{u}'_2 \neq \mathbf{u}_2$ (errors of types 3–5), let us define the pairwise error event

$$\begin{aligned} & \mathcal{A}(\mathbf{u}_1, \mathbf{u}'_1, \mathbf{u}_2, \mathbf{u}'_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) \\ &= [\mathcal{T}(Q_1) \times \mathcal{T}(Q_2)] \cap \\ & \quad \{(\mathbf{x}'_1, \mathbf{x}'_2) : P(\mathbf{u}'_1, \mathbf{u}'_2)W_1(\mathbf{y}_1|\mathbf{x}'_1)W_2(\mathbf{y}_2|\mathbf{x}'_2) \geq P(\mathbf{u}_1, \mathbf{u}_2)W_1(\mathbf{y}_1|\mathbf{x}_1)W_2(\mathbf{y}_2|\mathbf{x}_2)\}. \end{aligned}$$

The conditional error event, given $(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2, \mathcal{C}_1, \mathcal{C}_2)$, is given by

$$\begin{aligned} & \mathcal{E}(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2, \mathcal{C}_1, \mathcal{C}_2) \\ &= \bigcup_{\mathbf{u}'_1 \neq \mathbf{u}_1, \mathbf{u}'_2 \neq \mathbf{u}_2} \{P(\mathbf{u}'_1, \mathbf{u}'_2)W_1(\mathbf{y}_1|\mathbf{x}_1[\mathbf{u}'_1])W_2(\mathbf{y}_2|\mathbf{x}_2[\mathbf{u}'_2]) \geq \\ & \quad P(\mathbf{u}_1, \mathbf{u}_2)W_1(\mathbf{y}_1|\mathbf{x}_1[\mathbf{u}_1])W_2(\mathbf{y}_2|\mathbf{x}_2[\mathbf{u}_2])\} \\ &\triangleq \bigcup_{\mathbf{u}'_1 \neq \mathbf{u}_1, \mathbf{u}'_2 \neq \mathbf{u}_2} \mathcal{E}(\mathbf{u}_1, \mathbf{u}'_1, \mathbf{u}_2, \mathbf{u}'_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2, \mathcal{C}_1, \mathcal{C}_2) \end{aligned} \quad (76)$$

Here too, the exponential tightness of the truncated union bound for two dimensional unions of events with independence structure as above can be established using de Caen's lower bound [5] (see [29]). For errors of types 3 and 4, let us define

$$\begin{aligned} & \mathcal{A}_1(\mathbf{u}_1, \mathbf{u}'_1, \mathbf{u}_2, \mathbf{u}'_2, \mathbf{x}_1, \mathbf{y}_1) \\ &= \mathcal{T}(Q_1) \cap \{\mathbf{x}'_1 : P(\mathbf{u}'_1, \mathbf{u}'_2)W_1(\mathbf{y}_1|\mathbf{x}'_1) \geq P(\mathbf{u}_1, \mathbf{u}_2)W_1(\mathbf{y}_1|\mathbf{x}_1)\} \end{aligned} \quad (77)$$

and

$$\mathcal{A}_2(\mathbf{u}_1, \mathbf{u}'_1, \mathbf{u}_2, \mathbf{u}'_2, \mathbf{x}_2, \mathbf{y}_2)$$

$$= \mathcal{T}(Q_2) \cap \{x'_2 : P(u'_1, u'_2)W_2(y_2|x'_2) \geq P(u_1, u_2)W_2(y_2|x_2)\}. \quad (78)$$

The probability of $\mathcal{E}(u_1, u'_1, u_2, u'_2, x_1, x_2, y_1, y_2, \mathcal{C}_1, \mathcal{C}_2)$ (w.r.t. the randomness of the bin assignment) is given by:

$$\begin{aligned} & \overline{\Pr}\{\mathcal{E}(u_1, u'_1, u_2, u'_2, x_1, x_2, y_1, y_2, \mathcal{C}_1, \mathcal{C}_2)\} \\ &= 2^{-n(R_1+R_2)} \left| [\mathcal{C}_1 \times \mathcal{C}_2] \cap \mathcal{A}(u_1, u'_1, u_2, u'_2, x_1, x_2, y_1, y_2) \right| + \\ & \quad 2^{-n(R_1+R_2)} \left| \mathcal{C}_1 \cap \mathcal{A}_1(u_1, u'_1, u_2, u'_2, x_1, y_1) \right| + \\ & \quad 2^{-n(R_1+R_2)} \left| \mathcal{C}_2 \cap \mathcal{A}_2(u_1, u'_1, u_2, u'_2, x_2, y_2) \right| + \\ & \quad + 2^{-n(R_1+R_2)} \mathcal{I}\{P(u'_1, u'_2) \geq P(u_1, u_2)\}, \end{aligned} \quad (79)$$

where the first term stands for errors of type 5, the second and third terms represent errors of types 3 and 4, and the last term is associated with an error where both $u'_1 \neq u_1$ and $u'_2 \neq u_1$, but the respective bins both coincide. Passing temporarily to shorthand notation, let us denote

$$N \triangleq \left| [\mathcal{C}_1 \times \mathcal{C}_2] \cap \mathcal{A} \right| + \left| \mathcal{C}_1 \cap \mathcal{A}_1 \right| + \left| \mathcal{C}_2 \cap \mathcal{A}_2 \right| + \mathcal{I}\{P(u'_1, u'_2) \geq P(u_1, u_2)\} \triangleq N_{12} + N_1 + N_2 + I. \quad (80)$$

The next step, as before, is to average over the randomness of all codewords in \mathcal{C}_1 and \mathcal{C}_2 . To analyze the large deviations behavior of $N_{12} + N_1 + N_2$, the contributions of the individual random variables can be handled separately, since $\Pr\{N_{12} + N_1 + N_2 > \text{threshold}\}$ is of the same exponential order of the sum

$$\Pr\{N_{12} > \text{threshold}\} + \Pr\{N_1 > \text{threshold}\} + \Pr\{N_2 > \text{threshold}\}.$$

Now, N_1 and N_2 are binomial random variables whose numbers of trials are 2^{nR_1} and 2^{nR_2} , respectively, and whose probabilities of success decay exponentially according to the relevant channel mutual informations, similarly as before. So their contributions are again analyzed with great similarity to those of type 1 and type 2 errors.

Finally, it remains to handle N_{12} , which is not a binomial random variable, but it can be decomposed as the sum (over combinations of conditional types of x'_1 given y_1 and of x'_2 given y_2) of products of independent binomial random variables, for which we reuse the notations N_1 and N_2 (for a given combination of such types). Using the same techniques as in [21, Chap. 6], one can easily obtain the following generic result concerning the large deviations behavior of $N_1 \cdot N_2$: If N_1 is a binomial random variable with 2^{nA_1} trials and probability of success 2^{-nB_1} and N_2 is an independent binomial random variable with 2^{nA_2} trials and probability of success 2^{-nB_2} , then

$$\begin{aligned} \Pr\{N_1 \cdot N_2 \geq 2^{nC}\} & \doteq \max_{0 \leq \alpha \leq C} \Pr\{N_1 \geq 2^{n\alpha}\} \cdot \Pr\{N_2 \geq 2^{n(C-\alpha)}\} \\ & \doteq 2^{-nE} \end{aligned} \quad (81)$$

with

$$E = \begin{cases} [B_1 - A_1]_+ + [B_2 - A_2]_+ & C \leq [A_1 - B_1]_+ + [A_2 - B_2]_+ \\ \infty & C > [A_1 - B_1]_+ + [A_2 - B_2]_+ \end{cases} \quad (82)$$

Using this fact, it is possible to obtain the contribution of the type 5 error event.

Upon carrying out the analysis along these lines, the state of affairs turns out to be as described next. In the analysis of the conditional probability of error, the contribution of a given type class, $\mathcal{T}(\mathbf{u}'_1, \mathbf{u}'_2)$, of competing source vectors, which are encoded into \mathbf{x}'_1 and \mathbf{x}'_2 (from given conditional type classes given \mathbf{y}_1 and \mathbf{y}_2 , respectively) is the following: the probability of error of type i is of the exponential order of $\exp\{-n[f_i(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+\}$, $i = 1, \dots, 5$. Thus, the total conditional error probability contributed by this combination of types is of the exponential order of

$$\begin{aligned} \sum_{i=1}^5 \exp\{-n[f_i(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+\} & \doteq \exp\{-n \min_i [f_i(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+\} \\ & = \exp\{-n[f_0(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+\}. \end{aligned} \quad (83)$$

For the total contribution of all type classes, the exponent $[f_0(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+$ should be minimized over all such combinations of types (that yield the relevant pairwise error event). An upper bound on this exponent is obtained by selecting the same combination of types as those of the correct source vectors (instead of taking this minimum), namely, the conditional error probability of the MAP decoder is simply lower bounded by the exponential order of $\exp\{-n[f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)]_+\}$. As for the universal decoder, one should minimize the exponent $[f_0(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2)]_+$ as well, but only over the combinations of type classes that are associated with the pairwise error event of this decoder, namely, those for which $f_0(\mathbf{u}'_1, \mathbf{u}'_2, \mathbf{x}'_1, \mathbf{x}'_2, \mathbf{y}_1, \mathbf{y}_2) \geq f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)$. However, this minimum is exactly $[f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)]_+$, which agrees with that of the upper bound associated with the MAP decoder.

More formally, denoting $f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)$ as a functional of the relevant joint empirical distributions, i.e.,

$$F(\hat{P}_{\mathbf{u}_1 \mathbf{u}_2}, \hat{P}_{\mathbf{x}_1 \mathbf{y}_1}, \hat{P}_{\mathbf{x}_2 \mathbf{y}_2}) \triangleq [f_0(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2)]_+, \quad (84)$$

the error exponent achieved by both the MAP decoder and the universal decoder is given by

$$\begin{aligned} E(R_1, R_2, Q_1, Q_2) &= \min_{P_{U'_1 U'_2}, W'_1, W'_2} \left\{ D(P_{U'_1 U'_2} \| P_{U_1 U_2}) + D(W'_1 \| W_1 | Q_1) + \right. \\ &\quad \left. D(W'_2 \| W_2 | Q_2) + F(P_{U'_1 U'_2}, Q_1 \times W'_1, Q_2 \times W'_2) \right\}. \end{aligned} \quad (85)$$

Note that when R_1 and R_2 are sufficiently large, neither f_3 nor f_4 would achieve f_0 . At the same time, f_1 , f_2 and f_5 degenerate as follows:

$$f_1(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = \hat{I}(X_1; Y_1) - \hat{H}(U_1 | U_2) \quad (86)$$

$$f_2(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = \hat{I}(X_2; Y_2) - \hat{H}(U_2 | U_1) \quad (87)$$

$$f_5(\mathbf{u}_1, \mathbf{u}_2, \mathbf{x}_1, \mathbf{x}_2, \mathbf{y}_1, \mathbf{y}_2) = [\hat{I}(X_1; Y_1) + \hat{I}(X_2; Y_2) - \hat{H}(U_1, U_2)]_+. \quad (88)$$

Therefore, we have

$$\begin{aligned} E(\infty, \infty, Q_1, Q_2) &= \min_{P_{U'_1 U'_2}, W'_1, W'_2} \left[D(P_{U'_1 U'_2} \| P_{U_1 U_2}) + D(W'_1 \| W_1 | Q_1) + D(W'_2 \| W_2 | Q_2) + \right. \\ &\quad \min\{I(X_1; Y'_1) - H(U'_1 | U'_2), I(X_2; Y'_2) - H(U'_2 | U'_1), \\ &\quad \left. I(X_1; Y'_1) + I(X_2; Y'_2) - H(U'_1, U'_2)\} \right]. \end{aligned} \quad (89)$$

Further restricting this to the case of noiseless bit-pipes at fixed transmission rates r_1 and r_2 , respectively, the above channel-related divergence terms vanish, and one obtains the error exponent of separate compression and joint decompression of correlated sources

$$\min_{P_{U'_1 U'_2}} \left[D(P_{U'_1 U'_2} \| P_{U_1 U_2}) + \min\{r_1 - H(U'_1 | U'_2), r_2 - H(U'_2 | U'_1), r_1 + r_2 - H(U'_1, U'_2)\} \right], \quad (90)$$

in agreement with [4, Exercise 13.5] (second edition).

Appendix – Proof of Eq. (6)

First, observe that

$$(1 + \rho) \ln \left(\sum_{u \in \mathcal{U}} [P(u)]^{1/(1+\rho)} \right) = \max_{P'} [\rho \mathcal{H}(P') - D(P' \| P)], \quad (\text{A.1})$$

as can easily be seen by solving explicitly the maximization on the r.h.s. Next observe that for $\rho > 0$, the maximizer P'_0 is always associated with an entropy larger¹⁵ than $\mathcal{H}(P) = H(U)$. Therefore, the above identity can be further developed to obtain

$$\begin{aligned} &(1 + \rho) \ln \left(\sum_{u \in \mathcal{U}} [P(u)]^{1/(1+\rho)} \right) \\ &= \max_{\{P': \mathcal{H}(P) \leq \mathcal{H}(P') \leq \log |\mathcal{U}|\}} [\rho \mathcal{H}(P') - D(P' \| P)] \\ &= \max_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} \max_{\{P': \mathcal{H}(P') = \tilde{R}\}} [\rho \mathcal{H}(P') - D(P' \| P)] \\ &= \max_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} \max_{\{P': \mathcal{H}(P') = \tilde{R}\}} [\rho \tilde{R} - D(P' \| P)] \\ &= \max_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} \left[\rho \tilde{R} - \min_{\{P': \mathcal{H}(P') = \tilde{R}\}} D(P' \| P) \right] \\ &= \max_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} [\rho \tilde{R} - E^s(\tilde{R})], \end{aligned} \quad (\text{A.2})$$

¹⁵To see why this is true, note that $\rho \mathcal{H}(P) \leq \max_{P'} [\rho \mathcal{H}(P') - D(P' \| P)] = \rho \mathcal{H}(P'_0) - D(P'_0 \| P) \leq \rho \mathcal{H}(P'_0)$.

where the last step follows from the fact that, due to the convexity and the monotonicity of the source coding exponent function, the constraint $\mathcal{H}(P') \geq \tilde{R}$, of the minimization of $D(P' \| P)$ that defines it, is attained with equality in the range $\tilde{R} \in [\mathcal{H}(P), \log |\mathcal{U}|]$ (see also [2, eq. (7)]). Therefore,

$$\begin{aligned}
& \max_{0 \leq \rho \leq 1} \left\{ E_0(\rho, Q) - (1 + \rho) \ln \left(\sum_{u \in \mathcal{U}} [P(u)]^{1/(1+\rho)} \right) \right\} \\
&= \max_{0 \leq \rho \leq 1} \left[E_0(\rho, Q) - \max_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} [\rho \tilde{R} - E^s(\tilde{R})] \right] \\
&= \max_{0 \leq \rho \leq 1} \min_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} [E_0(\rho, Q) - \rho \tilde{R} + E^s(\tilde{R})] \\
&= \min_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} \max_{0 \leq \rho \leq 1} [E_0(\rho, Q) - \rho \tilde{R} + E^s(\tilde{R})] \\
&= \min_{\mathcal{H}(P) \leq \tilde{R} \leq \log |\mathcal{U}|} [E_r^c(\tilde{R}, Q) + E^s(\tilde{R})], \tag{A.3}
\end{aligned}$$

where the interchange of the maximization and the minimization in the second to the last step is allowed by the concavity of $E_0(\rho, Q)$ in ρ [9, eq. (5.6.26)] and the convexity of $E^s(\tilde{R})$ [2].

Acknowledgement

Interesting discussions with Jacob Ziv are acknowledged with thanks.

References

- [1] J. Chen, D.-k. He, A. Jagmohan, and L. A. Lastras-Montaño, “On universal variable-rate Slepian–Wolf coding,” *Proc. 2008 IEEE International Conference on Communications (ICC 2008)*, pp. 1426–1430, 2008.
- [2] I. Csiszár, “Joint source–channel error exponent,” *Problems of Control and Information Theory*, vol. 9, no. 5, pp. 315–328, 1980.
- [3] I. Csiszár, “Linear codes for sources and source networks: error exponents, universal coding,” *IEEE Trans. Inform. Theory*, vol. IT–28, no. 4, pp. 585–592, July 1982.
- [4] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Academic Press, 1981. Second Edition: Cambridge University Press, New York, 2011.
- [5] D. de Caen, “A lower bound on the probability of a union,” *Discrete Math.*, vol. 169, pp. 217–220, 1997.
- [6] S. C. Draper, “Universal incremental Slepian–Wolf coding,” *Proc. 42nd Annual Allerton Conference on Communication, Control and Computing*, Monticello, IL, USA, October 2004.
- [7] M. Feder and A. Lapidoth, “Universal decoding for channels with memory,” *IEEE Trans. Inform. Theory*, vol. 44, no. 5, pp. 1726–1745, September 1998.

- [8] M. Feder and N. Merhav, “Universal composite hypothesis testing: a competitive minimax approach,” *IEEE Trans. Inform. Theory*, special issue in memory of Aaron D. Wyner, vol. 48, no. 6, pp. 1504–1517, June 2002.
- [9] R. G. Gallager, *Information Theory and Reliable Communication*, John Wiley & Sons, New York 1968.
- [10] V. D. Goppa, “Nonprobabilistic mutual information without memory,” *Probl. Cont. Information Theory*, vol. 4, pp. 97–102, 1975.
- [11] F. Jelinek, *Probabilistic Information Theory*, McGraw Hill, New York 1968.
- [12] J. C. Kieffer, “Some universal noiseless multiterminal source coding theorems,” *Information and Control*, vol. 46, pp. 93–107, 1980.
- [13] A. Lapidoth and J. Ziv, “On the universality of the LZ-based noisy channels decoding algorithm,” *IEEE Trans. Inform. Theory*, vol. 44, no. 5, pp. 1746–1755, September 1998.
- [14] Y. Lomnitz and M. Feder, “Communication over individual channels – a general framework,” arXiv:1023.1406v1 [cs.IT] 7 Mar 2012.
- [15] Y. Lomnitz and M. Feder, “Universal communication over modulo-additive channels with an individual noise sequence,” arXiv:1012.2751v2 [cs.IT] 7 May 2012.
- [16] K. Marton, “Error exponent for source coding with a fidelity criterion,” *IEEE Trans. Inform. Theory*, vol. IT-20, no. 2, pp. 197–199, March 1974.
- [17] S. Matloub and T. Weissman, “Universal zero-delay joint source-channel coding,” *IEEE Trans. Inform. Theory*, vol. 52, no. 12, pp. 5240–5250, December 2006.
- [18] N. Merhav, “Universal decoding for memoryless Gaussian channels with a deterministic interference,” *IEEE Trans. Inform. Theory*, vol. 39, no. 4, pp. 1261–1269, July 1993.
- [19] N. Merhav, “Universal detection of messages via finite-state channels,” *IEEE Trans. Inform. Theory*, vol. 46, no. 6, pp. 2242–2246, September 2000.
- [20] N. Merhav, “Shannon’s secrecy system with informed receivers and its application to systematic coding for wiretapped channels,” *IEEE Trans. Inform. Theory*, special issue on *Information-Theoretic Security*, vol. 54, no. 6, pp. 2723–2734, June 2008.
- [21] N. Merhav, “Statistical physics and information theory,” *Foundations and Trends in Communications and Information Theory*, vol. 6, nos. 1–2, pp. 1–212, 2009.
- [22] N. Merhav, “Relations between random coding exponents and the statistical physics of random codes,” *IEEE Trans. Inform. Theory*, vol. 55, no. 1, pp. 83–92, January 2009.
- [23] N. Merhav, “Erasure/list exponents for Slepian–Wolf decoding,” *IEEE Trans. Inform. Theory*, vol. 60, no. 8, pp. 4463–4471, August 2014.
- [24] N. Merhav, “Exact random coding exponents of optimal bin index decoding,” *IEEE Trans. Inform. Theory*, vol. 60, no. 10, pp. 6024–6031, October 2014.

- [25] N. Merhav, “Universal decoding for arbitrary channels relative to a given family of decoding metrics,” *IEEE Trans. Inform. Theory*, vol. 59, no. 9, pp. 5566–5576, September 2013.
- [26] V. Misra and T. Weissman, “The porosity of additive noise sequences,” arXiv:1025.6974v1 [cs.IT] 31 May 2012.
- [27] Y. Oohama and T. S. Han, “Universal coding for the Slepian–Wolf data compression system and the strong converse theorem,” *IEEE Trans. Inform. Theory*, vol. 40, no. 6, pp. 1908–1919, November 1994.
- [28] S. Sarvotham, D. Baron, and R. G. Baraniuk, “Variable–rate universal Slepian–Wolf coding with feedback,” *Proc. 39th Asilomar Conference on Signals, Systems and Computers*, pp. 8–12, November 2005.
- [29] J. Scarlett, A. Martínéz, and A. i. Fábregas, “Multiuser techniques for mismatched decoding,” submitted to *IEEE Trans. Inform. Theory*, November 2013. arxiv.org/pdf/1311.6635
- [30] S. Shamai (Shitz), S. Verdú and R. Zamir, “Systematic lossy source/ channel coding,” *IEEE Trans. Inform. Theory*, vol. 44, no. 2, pp. 564–579, March 1998.
- [31] N. Shulman, *Communication over an Unknown Channel via Common Broadcasting*, Ph.D. dissertation, Department of Electrical Engineering – Systems, Tel Aviv University, July 2003.
- [32] N. Weinberger and N. Merhav, “Optimum trade-off between the error exponent and the excess–rate exponent of variable–rate Slepian–Wolf coding,” *IEEE Trans. Inform. Theory*, vol. 61, no. 4, pp. 2165–2190, April 2015.
- [33] J. Ziv, “Universal decoding for finite–state channels,” *IEEE Trans. Inform. Theory*, vol. IT–31, no. 4, pp. 453–460, July 1985.