

On Successive Refinement for the Wyner–Ziv Problem

Yossef Steinberg and Neri Merhav

Department of Electrical Engineering
Technion - Israel Institute of Technology
Haifa 32000, ISRAEL
[ysteinbe,merhav]@ee.technion.ac.il

March 19, 2003

Abstract

Achievable rates are characterized for successive refinement in the Wyner–Ziv scenario, namely, in the presence of correlated side information (SI) at the receiver. In this setting, the encoder is assumed to operate in two stages, where the first corresponds to relatively low rate and high distortion, and the second, comprising a refinement code on top of the first code, is aimed at reproduction at reduced distortion. Both decoders (for low-rate/high-distortion and for high-rate/low-distortion) are equipped with SI streams, correlated to the source, but unavailable to the encoder. Furthermore, it is assumed that the decoder that receives the higher rate bitstream, i.e., the additional refinement bits, accesses also SI of better quality (in a sense that will be defined later) than that of the lower resolution decoder. For a memoryless joint process (that includes the source to be encoded and its instantaneously correlated SI streams), necessary and sufficient conditions are furnished, in terms of single-letter formulas, for the achievability of a pair of rates, corresponding to two given distortion levels. Special attention is devoted to the degenerate, but important, case where the two SI streams, at the two decoders, are identical. For this case, conditions are provided for successive refinability in the sense of the existence of codes that asymptotically achieve the Wyner–Ziv rate–distortion function, simultaneously at both distortion levels. In this context, the doubly symmetric binary source (with the Hamming distortion measure) and the jointly Gaussian source (with the squared error distortion measure) are shown to be successively refinable in the Wyner–Ziv setting. It is also demonstrated that a source that is not successively refinable in the ordinary sense (i.e., without SI) may become successively refinable in the presence of SI at the decoders.

Index terms — Side information, successive refinement, scalable coding, progressive coding, multiple description, Wyner-Ziv problem.

1 Introduction

The problem of successive refinement of information was originally formulated by Koshelev [7], and by Equitz and Cover [4], who viewed this problem as a special case of the more general multiple description problem. In their original setting, a source \mathbf{X} is to be encoded and transmitted to its destination over a rate-limited channel. To this end, due to the fact that the channel is rate-limited, the encoder produces a compressed bit string \mathbf{S} at rate R_1 , from which the decoder produces $\hat{\mathbf{X}}$, an approximation of the original source at distortion level D_1 , according to some distortion measure $d(\cdot, \cdot)$. Rate-distortion theory tells us, as is well known, that the minimal rate R_1 that is needed in order to enable reproduction of the source at accuracy D_1 is given by $R_X(D_1)$, the rate-distortion function of the source at distortion level D_1 . At a later stage, a more accurate description of the source is needed, and the encoder sends a secondary string of compressed bits, \mathbf{S}_Δ , at rate ΔR , to the destination. The decoder, having at hand both strings \mathbf{S} and \mathbf{S}_Δ , produces a more accurate reproduction of the source, $\tilde{\mathbf{X}}$, at distortion level D_2 ($D_2 < D_1$). Invoking again the fundamental limits of rate-distortion theory, the total rate $R_1 + \Delta R$ cannot be lower than $R_X(D_2)$, the rate-distortion function evaluated at the corresponding distortion level. The best one can hope for, then, is that the two rates *simultaneously* lie on the rate-distortion curve, i.e.,

$$R_1 = R_X(D_1), \quad \text{and} \quad R_1 + \Delta R = R_X(D_2). \quad (1)$$

This, however, is not always possible. In general, some penalty might have to be paid for successive coding. In the ideal situation, where successive coding, in two or more stages, can be made rate-distortion optimal simultaneously at all stages, the source is called *successively refinable*. Koshelev [7], [8], and Equitz and Cover [4], have shown that a necessary and sufficient condition for a source to be successively refinable is that the conditional distributions $P_{\tilde{X}|X}$ and $P_{\hat{X}|X}$ that achieve the rate distortion function at distortion levels D_2 and D_1 , respectively, are Markov compatible in the sense that they can be represented as a Markov chain $X \ominus \tilde{X} \ominus \hat{X}$. Equitz and Cover have shown a few examples of successively refinable sources. In particular, it was shown that Gaussian sources with the squared error distortion measure, general discrete memoryless sources (DMS's) with the Hamming distortion measure, and Laplacian sources with the absolute distortion measure, are all successively refinable. Moreover, an example of a source that does not satisfy the Markov condition was given, thus showing that the problem is not redundant, i.e., not every source is successively refinable.

Clearly, the theoretical and practical value of successive coding is not limited to successively refinable sources. Successive coding is suitable to any application where a relatively coarse description of the source suffices at the first use of the data, and fine details are needed only at some later stage. Transmission of the fine details, which comprises the refinement stage, may take place upon user request or upon availability of additional system resources, like free time slots or extra bandwidth in a network communication environment, reduction of cost of channel use, etc. In [7],¹ Koshelev characterized the set of all quadruples $(R_1, R_1 + \Delta R, D_1, D_2)$ that are achievable via successive coding for a general DMS (that is not necessarily successively refinable). From this more general characterization, the Markov condition for successive refinability is readily obtained as a byproduct.

The above description of earlier work on successive refinement of information serves as one part of the background for this work. The other part has to do with rate–distortion coding with side information (SI) at the decoder, which is well–known as the Wyner–Ziv problem [12], the lossy counterpart of a certain version of the Slepian–Wolf setting. In some source coding applications, it is plausible to assume that the decoder has some knowledge about the encoded source, which is not available to the encoder. For example, consider an airplane taking aerial photographs of a certain area, for the purpose of later processing (e.g., updating maps) on the ground. The photographs are encoded and transmitted to a base-station, and on the ground, the decoder reconstructs the aerial photos. For that purpose, the decoder takes advantage of previously existing data about the photographed area – for example, previous photographs, maps, etc. Since the equipment on the plane should be light and compact, the encoder is kept as simple as possible, and thus it does not utilize previous photographs or maps of the area in the encoding process. Another application of source coding with SI at the decoder is that of systematic source/channel coding, suggested by Shamai and Verdú for lossless coding [10], and extended later to lossy coding in [11]. In their setting, an analog information source is to be transmitted to the receiver and reconstructed at some distortion level D . For that purpose, the sender has access to two channels: an analog channel, over which the source is transmitted uncodedly (the systematic part), and a digital channel, over which digital, coded information about the source is sent. The receiver has access to the two outputs. The output of the analog channel is the original analog source corrupted by noise. By contrast, essentially no errors are introduced by the digital channel, as the information is channel-coded. Such a model arises, for example, when an analog communication system is upgraded with a digital

¹See also Rimoldi [9].

system, but back-compatibility for users utilizing only the analog part is being kept. Now, observe that the output of the analog channel serves as SI about the source, which is accessible to the decoder but not to the encoder.

The main result of Wyner and Ziv [12] is as follows: Let \mathbf{Z} be the SI, available to the decoder only, and having joint distribution P_{XZ} with the encoded source \mathbf{X} . The encoder, although ignorant of the specific realization of the SI, has full knowledge of this joint distribution. Let $R_{X|Z}^*(D)$ stand for the minimal achievable rate for coding with SI. Wyner and Ziv have shown that

$$R_{X|Z}^*(D) = \min I(X; U|Z), \quad (2)$$

where the minimization is over all random variables U such that $U \circlearrowleft X \circlearrowleft Z$ is a Markov chain and such that there exists a function f of U and Z such that

$$\mathbb{E}d(X, f(U, Z)) \leq D. \quad (3)$$

This work is a joint extension of the results of [7], [8], and [9] on the one hand, and the results of [12] on the other hand. Analogously to the previous works on successive coding, where classical rate-distortion theory forms the base of reference, here this role is played by the Wyner–Ziv result (2), (3). In particular, let R_1 , ΔR , and D_1 , D_2 be the rates and distortion levels of the coarse and refinement stages, respectively, as defined for the classical successive refinement above. Let \mathbf{Z} , \mathbf{Y} be the SI available to the decoder at the coarse and at the refinement stages, respectively (see Fig. 1). The natural question that may be asked, in this setting, is what is the set of all achievable quadruples $(R_1, \Delta R, D_1, D_2)$ and how to achieve every such quadruple. Note that for a general joint source, $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, there is no guarantee that $D_2 \leq D_1$ since the SI available at the first stage, \mathbf{Z} , may be better than the SI available at the second stage, \mathbf{Y} . Indeed, this seems to be quite a difficult problem, and we have not been able to fully characterize the achievable rates and distortion levels when no structure is imposed on the joint distribution of the source and SI. However, we have been able to fully characterize the set of all achievable quadruples $(R_1, \Delta R, D_1, D_2)$ for the special case where the Markov relation $\mathbf{X} \circlearrowleft \mathbf{Y} \circlearrowleft \mathbf{Z}$ hold, namely, the SI \mathbf{Y} can be considered to be of better quality than \mathbf{Z} . From the practical point of view, this structure is not too restrictive. To illustrate this point, let \mathbf{Z} , \mathbf{Z}_Δ be the SI *received* at the destination at the coarse and at the refinement stages, respectively, and assume that for the purpose of decoding at the refinement stage, the decoder has at hand both \mathbf{Z} and \mathbf{Z}_Δ , i.e., $\mathbf{Y} = (\mathbf{Z}, \mathbf{Z}_\Delta)$. Then the Markov structure $\mathbf{X} \circlearrowleft \mathbf{Y} \circlearrowleft \mathbf{Z}$ does hold.

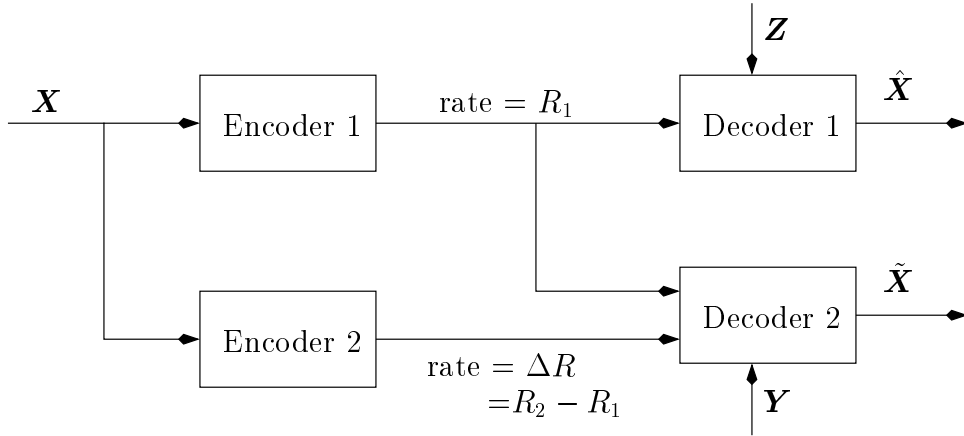


Figure 1: Successive refinement with side information (SI) at the decoder.

As an application of successive coding with increasing levels of SI, consider again the setup of [11]: We have an analog communication channel, which was upgraded with a digital channel. Since back-compatibility for subscribers utilizing only the analog part is kept, coded source-words are sent over the digital channel and, simultaneously, the source is sent uncodedly over the analog channel. The output of the analog channel serves as SI at the decoder. In this setup, whenever refinement information is sent over the digital channel (by user request, or by availability of digital channel resources), additional SI \mathbf{Z}_Δ arrives as well at the decoder via the analog channel. If all past SI is being stored at the destination, we have a Wyner–Ziv successive coding scheme with the Markov structure $\mathbf{X} \leftrightarrow \mathbf{Y} \leftrightarrow \mathbf{Z}$.

Analogously to (1), a source is successively refinable if

$$R_1 = R_{X|Z}^*(D_1), \quad \text{and} \quad R_1 + \Delta R = R_{X|Y}^*(D_2). \quad (4)$$

Note, however, that the codewords sent at the first stage are supposed to be decoded with the “weak” SI \mathbf{Z} , and thus some rate might be lost relative to the case where the vector \mathbf{Y} is available to the decoder already at the first stage. In such case, the source cannot be successively refinable. We show later that a necessary condition for successive refinability to hold, is that the SI \mathbf{Y} is *equivalent* (in a sense that will be made clear in Section 3.2) to the SI \mathbf{Z} at distortion level D_1 . One special case of such equivalence is when the SI random vectors are identical, i.e., $\mathbf{Y} = \mathbf{Z}$. However, there are also other cases where successive refinability holds, and the concept of equivalence of SI turns out to be more general than the requirement $\mathbf{Y} = \mathbf{Z}$.

A problem related to successive coding with SI is that of coding when SI may be present (or absent), which was studied by Heegard and Berger [5] and independently by Kaspi [6]. Specifically, suppose that in the classical Wyner–Ziv setting (i.e., no successive refinement), two levels of SI may be available to the decoder: \mathbf{Z} , or a “better” description \mathbf{Y} . The encoder strives to guarantee two levels of distortion: D_1 if only \mathbf{Z} is available, and $D_2 < D_1$ if the better SI \mathbf{Y} is available. Heegard and Berger [5] derived a single-letter formula for the minimal rate needed to guarantee the pair of distortion levels (D_1, D_2) . Clearly, this model is a special case of the model considered here with $R_2 = R_1$ (cf. Fig. 1), and it should be pointed out that even in this case, Heegard and Berger could find full characterization of achievable performance only under the assumption of a Markov structure $X \leftrightarrow Y \leftrightarrow Z$, and the general case, of a joint distribution of (X, Y, Z) , was left open. A–fortiori, it would then be hard to expect a complete characterization for $R_2 \geq R_1$.

This paper is organized as follows: Notation and definitions are provided in Section 2. The main results are presented and discussed in Section 3. Examples of sources which are successively refinable with SI are given in Section 4. In Section 5, we present an extension of part of the results of Section 3 to multistage successive coding schemes. Finally, the proofs of the main theorems are given in Section 6.

2 Notation and Preliminaries

We begin by setting up the notation. Let \mathcal{X} be a finite set and let \mathcal{X}^n be the set of all n -vectors with components in \mathcal{X} . A member of \mathcal{X}^n will be written as $x^n = (x_1, x_2, \dots, x_n)$, and substrings of x^n are written as $x_i^j = (x_i, x_{i+1}, \dots, x_j)$ for $i \leq j$. When $i = 1$, the subscript will be omitted. When the dimension is clear from the context, vectors will be denoted by boldface letters, e.g., $\mathbf{x} \in \mathcal{X}^n$. A similar convention is used for random variables and vectors, which are denoted by upper case letters. A discrete memoryless source (DMS) (\mathcal{X}, P_X) is an infinite sequence $\{X_i\}_{i=1}^\infty$ of independent copies of a random variable X taking values in \mathcal{X} with a generic distribution P_X , i.e.,

$$P_X(x^n) = \prod_{i=1}^n P_X(x_i).$$

Similarly, a triple source $(\mathcal{X}\mathcal{Y}\mathcal{Z}, P_{XYZ})$ is an infinite sequence of independent copies of the triplet of random variables (X, Y, Z) taking values in the finite sets \mathcal{X} , \mathcal{Y} , and \mathcal{Z} , respectively, with generic joint distribution P_{XYZ} . The induced marginals and conditional distributions are denoted

by the corresponding subscripts, e.g., P_X , P_{YZ} , $P_{Z|X}$, etc. Whenever clear from the context, these subscripts will be omitted, e.g., $P_{XZ}(x, z)$ will sometimes be denoted simply by $P(x, z)$. Also, with a slight abuse of notation, and when there is no room for ambiguity, we will denote a source (\mathcal{X}, P_X) by referring to its generic distribution P_X or random variable X .

We are interested in coding the source X . Let $\hat{\mathcal{X}}$ stand for a finite reconstruction alphabet, and let

$$d: \mathcal{X} \times \hat{\mathcal{X}} \rightarrow [0, \infty)$$

be a single-letter distortion measure. The vector distortion measure is defined, in the usual way, as

$$d(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i) \quad \forall \mathbf{x} \in \mathcal{X}^n, \hat{\mathbf{x}} \in \hat{\mathcal{X}}^n.$$

Definition 1 An (n, M_1, M_2, D_1, D_2) *successive refinement (SR) code* for the source X with SI (Y, Z) consists of a first-stage encoder–decoder pair (ϕ_1, ψ_1) :

$$\begin{aligned} \phi_1: \quad \mathcal{X}^n &\rightarrow \{1, 2, \dots, M_1\} \\ \psi_1: \quad \{1, 2, \dots, M_1\} \times \mathcal{Z}^n &\rightarrow \hat{\mathcal{X}}^n, \end{aligned} \tag{5}$$

and a second-stage (or *refinement*) encoder–decoder pair (ϕ_2, ψ_2) :

$$\begin{aligned} \phi_2: \quad \mathcal{X}^n &\rightarrow \{1, 2, \dots, M_2\} \\ \psi_2: \quad \{1, 2, \dots, M_1\} \times \{1, 2, \dots, M_2\} \times \mathcal{Y}^n &\rightarrow \hat{\mathcal{X}}^n, \end{aligned} \tag{6}$$

such that

$$\mathbb{E}d(X^n, \psi_1(\phi_1(X^n), Z^n)) \leq D_1 \tag{7}$$

and

$$\mathbb{E}d(X^n, \psi_2(\phi_1(X^n), \phi_2(X^n), Y^n)) \leq D_2, \tag{8}$$

where \mathbb{E} stands for the expectation operation.

Remark 1 One may consider a slightly more general definition, without any difficulty in the analysis, where different distortion measures are used in the two stages, i.e., d_1 in (7) and d_2 in (8). For the sake of simplicity, however, we use the same distortion measure at both stages.

The *rate pair* (R_1, R_2) of the (n, M_1, M_2, D_1, D_2) SR code is

$$\begin{aligned} R_1 &= \frac{1}{n} \log M_1, \\ R_2 &= \frac{1}{n} \log(M_1 M_2). \end{aligned} \tag{9}$$

Given a distortion pair $\mathbf{D} = (D_1, D_2)$, a rate pair (R_1, R_2) is said to be \mathbf{D} -*achievable* for X with SI (Y, Z) if, for any $\delta > 0$, $\epsilon > 0$, and sufficiently large n , there exists an $(n, \exp[n(R_1 + \delta)], \exp[n(R_2 - R_1 + \delta)], D_1 + \epsilon, D_2 + \epsilon)$ SR code for the source X with SI (Y, Z) . The collection of all \mathbf{D} -achievable rate pairs is the achievable SR region for coding with SI, and is denoted by $\mathcal{R}(\mathbf{D})$.

In the sequel, special attention will be given to the case where the SI, available at the two stages, is identical, i.e., $\mathcal{Y} = \mathcal{Z}$ and $P_{XY} = P_{XZ}$. Therefore, for convenience, we will denote the achievable SR region for coding with identical SI in the two stages by $\mathcal{R}_i(\mathbf{D})$.

An immediate consequence of Definition 1 is that $\mathcal{R}(\mathbf{D})$ depends on the joint distribution P_{XYZ} only via the marginals P_{XY} and P_{XZ} . Following the terminology in [5] (and also common in the context of broadcast channels [1]), a source P_{XYZ} is said to be *stochastically degraded* if exists a source $P_{\tilde{X}\tilde{Y}\tilde{Z}}$ such that $\tilde{X} \rightarrow \tilde{Y} \rightarrow \tilde{Z}$ is a Markov chain, $P_{\tilde{X}\tilde{Y}} = P_{XY}$, and $P_{\tilde{X}\tilde{Z}} = P_{XZ}$. Since the achievable rate region $\mathcal{R}(\mathbf{D})$ depends only on the pair marginals, no distinction has to be made between physically degraded (Markov structured) sources and stochastically degraded sources. In particular, in the case of identical SI, that was mentioned in the previous paragraph, we might as well assume that $Y = Z$ with probability one, without loss of generality. For a general stochastically degraded source, we will occasionally use the terms *good* and *bad* SI to signify Y and Z , respectively.

3 Main Results

In this section, we present a single-letter characterization of the region of achievable rates for successive coding with Markov structured SI. Based on this characterization, we derive necessary and sufficient conditions for successive refinability.

3.1 The Achievable Region

Let a distortion pair $\mathbf{D} = (D_1, D_2)$ be given. Define $\mathcal{R}^*(\mathbf{D})$ to be the set of all rate pairs (R_1, R_2) for which there exists a triple of random variables (U, V, W) , taking values in finite alphabets, \mathcal{U} , \mathcal{V} , \mathcal{W} , respectively, such that the following conditions are satisfied:

1. $(U, V, W) \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$ is a Markov chain.

2. There exist deterministic maps

$$f_1 : \mathcal{U} \times \mathcal{Z} \rightarrow \hat{\mathcal{X}} \tag{10}$$

$$f_2 : \mathcal{W} \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}, \tag{11}$$

such that

$$\mathbb{E}d(X, f_1(U, Z)) \leq D_1 \tag{12}$$

$$\mathbb{E}d(X, f_2(W, Y)) \leq D_2 \tag{13}$$

3. The alphabets $\mathcal{U}, \mathcal{V}, \mathcal{W}$ satisfy:

$$|\mathcal{U}| \leq |\mathcal{X}| + 2, \tag{14}$$

$$|\mathcal{V}| \leq (|\mathcal{X}| + 1)^2, \tag{15}$$

$$|\mathcal{W}| \leq |\mathcal{X}|(|\mathcal{X}| + 2)(|\mathcal{X}| + 1)^2 + 1. \tag{16}$$

4. The rates R_1 and R_2 satisfy

$$R_1 \geq I(X; U|Z) + I(X; V|U, Y), \tag{17}$$

$$R_2 - R_1 \geq I(X; W|U, V, Y). \tag{18}$$

Our main result is the following:

Theorem 1 For any discrete memoryless stochastically degraded source, $X \circlearrowleft Y \circlearrowleft Z$,

$$\mathcal{R}(\mathbf{D}) = \mathcal{R}^*(\mathbf{D}).$$

The proof appears in Section 6.

Remark 2 We take this opportunity to point out an inaccuracy in the presentation of the main results of [7] and [9]. Koshelev and Rimoldi present the achievable rates as

$$\begin{aligned} R_1 &\geq I(X; \hat{X}_1), & \mathbb{E}d(X, \hat{X}_1) &\leq D_1, \\ R_2 &\geq I(X; \hat{X}_1 \hat{X}_2), & \mathbb{E}d(X, \hat{X}_2) &\leq D_2. \end{aligned}$$

Observe, however, that such a presentation includes in the achievable region rate pairs (R_1, R_2) for which $R_1 > R_2$. This is impossible by the very definition of our coding setup, as R_i , $i = 1, 2$ stand for cumulative rates, and thus loss of rate caused by poor design of the first stage, cannot be compensated for at the second, refinement stage. Note that in [9, Lemma 3], the direct (achievability) part is proved for R_1 and the differential rate $R_2 - R_1$, and has the form of our result with degenerate side information (Y, Z) . The upper bound, on the other hand, is proved for the cumulative rates [9, Lemma 1]. The crux of the problem lies in Lemma 4 there, where achievability of the cumulative rates is claimed, based on the (correct) result on differential rates.

Remark 3 Due to the conditioning on U in (17) and on U and V in (18), we can add structure to the auxiliary random variables in the above definition of $\mathcal{R}^*(\mathbf{D})$, at the expense of increased alphabet sizes. Specifically, define $\tilde{V} = UV$, $\tilde{W} = UVW$. Then we have, instead of condition 1 in the definition of $\mathcal{R}^*(\mathbf{D})$, the following structure

$$U \circlearrowleft \tilde{V} \circlearrowleft \tilde{W} \circlearrowleft X \circlearrowleft Y \circlearrowleft Z \tag{19}$$

and the bounds on the rates are now read

$$R_1 \geq I(X; U|Z) + I(X; \tilde{V}|U, Y), \tag{20}$$

$$R_2 - R_1 \geq I(X; \tilde{W}|\tilde{V}, Y). \tag{21}$$

It is easy to verify that the distortion levels are not altered by this substitution. Adding this Markov structure results in increased alphabet sizes $|\tilde{\mathcal{V}}| = |\mathcal{U}| \cdot |\mathcal{V}|$, and $|\tilde{\mathcal{W}}| = |\mathcal{U}| \cdot |\mathcal{V}| \cdot |\mathcal{W}|$.

Discussion: Observe that if $R_2 = R_1$, i.e., there is no excess rate in the refinement stage, $\mathcal{R}^*(\mathbf{D})$ coincides with the region given in [5, Theorem 3]. The random variables U and V in eqs. (17) and (18) represent the information about X that is sent to the decoder at the first stage, at rates $I(X; U|Z)$ and $I(X; V|U, Y)$, respectively. Due to the conditioning on Y , the codewords represented

by V can be decoded only at the refinement stage. For a given choice of (U, V, W) that satisfies the above conditions, let us write

$$R_1 = R_{1,1} + R_{1,2} \quad (22)$$

where

$$R_{1,1} = I(X; U|Z) \quad (23)$$

$$R_{1,2} = I(X; V|U, Y). \quad (24)$$

Thus, for given R_1 , there is a tradeoff between the part of the rate that can be utilized at the first stage, $R_{1,1}$, and the part that must wait to be decoded at the refinement stage, $R_{1,2}$. This tradeoff is controlled by the choice of the random variables U and V . One may question the wisdom of sending at the first stage information that can be decoded only at the refinement stage. Note, however, that coding rate can be better utilized at the second stage due to the presence of better SI, Y . Therefore, coding with $R_{1,2} > 0$ is beneficial whenever we want to utilize Y beyond what is possible with the limited refinement rate $R_2 - R_1$. Since V can be decoded only at the second stage, there is a degree of freedom in deciding whether to send it in the first or the second stage. This can be seen from eq. (17) as follows: Let (R_1, R_2, D_1, D_2) be achievable with a certain triplet (U, V, W) , and decompose R_1 as in (22)–(24). Define $\tilde{W} = VW$, let \tilde{V} be a null random variable (e.g., constant), and let $(\tilde{R}_1, \tilde{R}_2, \tilde{D}_1, \tilde{D}_2)$ be the rates and distortion levels achievable with the triplet $(U, \tilde{V}, \tilde{W})$. Since the functions f_1, f_2 do not depend on V , the distortion levels do not change, that is

$$(\tilde{D}_1, \tilde{D}_2) = (D_1, D_2) \quad (25)$$

and for the new rates, we have

$$\tilde{R}_1 = R_1 - R_{1,2} = R_{1,1} \quad (26)$$

$$\tilde{R}_2 = R_2, \quad (27)$$

$$\tilde{R}_2 - \tilde{R}_1 = R_2 - R_1 + R_{1,2}. \quad (28)$$

We now turn to the important special case of identical SI. Let us define $\mathcal{R}_i^*(\mathbf{D})$ similarly as $\mathcal{R}^*(\mathbf{D})$ where $Y \equiv Z$, the bounds on the alphabet sizes (14) and (16) are replaced by

$$|\mathcal{U}| \leq |\mathcal{X}| + 2 \quad (29)$$

$$|\mathcal{W}| \leq |\mathcal{X}|(|\mathcal{X}| + 2) + 1 \quad (30)$$

and the rate inequalities (17) and (18) are replaced by the following inequalities:

$$R_1 \geq I(X; U|Y), \quad (31)$$

$$R_2 - R_1 \geq I(X; W|UY). \quad (32)$$

Defining $\mathcal{R}_i(\mathbf{D})$ analogously to $\mathcal{R}(\mathbf{D})$ with the restriction of identical SI, we now have the following corollary to Theorem 1:

Corollary 1 For a discrete memoryless joint source (X, Y) :

$$\mathcal{R}_i(\mathbf{D}) = \mathcal{R}_i^*(\mathbf{D}).$$

Proof. In view of Theorem 1, we have to show that when $Y = Z$, the characterization of $\mathcal{R}^*(\mathbf{D})$ reduces to that of $\mathcal{R}_i^*(\mathbf{D})$. Indeed, if $Y = Z$, by (17), (18) we can write

$$R_1 \geq I(X; U|Z) + I(X; V|U, Y) = I(X; UV|Y) \quad (33)$$

$$R_2 - R_1 \geq I(X; W|UVY). \quad (34)$$

Observe that the auxiliary random variables U and V appear in the mutual information functions in (33), (34) always as a pair. The functions f_1 and f_2 do not depend on V . Thus we can define a new auxiliary random variable $\tilde{U} = UV$ and a new function $\tilde{f}_1(\tilde{U}) = f_1(U)$, without altering the distortions or the rates. This proves the corollary with the only exception that the alphabet size of \mathcal{W} is given by (16) instead of (30), and that of $\tilde{\mathcal{U}}$ is given by

$$|\tilde{\mathcal{U}}| = |\mathcal{U}| \cdot |\mathcal{V}| \quad (35)$$

instead of the right hand side of (29). The fact that we can restrict the alphabet sizes to those given by (29) and (30) follows from Carathéodory's theorem, in a manner similar to the bounds on the alphabet sizes in the proof of Theorem 1. We omit these details here. \square

3.2 Successive Refinability

The notion of successive refinability with SI is now defined.

Definition 2 A source X is said to be *successively refinable from D_1 to D_2 ($D_1 > D_2$) with SI* if

$$(R_{X|Z}^*(D_1), R_{X|Y}^*(D_2)) \in \mathcal{R}(D_1, D_2). \quad (36)$$

Similarly, a source X is said to be *successively refinable from D_1 to D_2 ($D_1 > D_2$) with identical SI Y* if

$$(R_{X|Y}^*(D_1), R_{X|Y}^*(D_2)) \in \mathcal{R}_i(D_1, D_2). \quad (37)$$

A characterization of sources which are successively refinable with SI, in the spirit of the results of Koshelev and Equitz and Cover, is given next.

Theorem 2 A source X with degraded SI (Y, Z) is successively refinable from D_1 to D_2 if and only if there exist a pair of random variables (U, W) and a pair of deterministic maps $f_1 : \mathcal{U} \times \mathcal{Z} \rightarrow \hat{\mathcal{X}}$, and $f_2 : \mathcal{W} \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}$, such that the following conditions simultaneously hold:

1. $R_{X|Z}^*(D_1) = I(X; U|Z)$ and $\mathbb{E}d(X, f_1(U, Z)) \leq D_1$,
2. $R_{X|Y}^*(D_2) = I(X; W|Y)$ and $\mathbb{E}d(X, f_2(W, Y)) \leq D_2$,
3. $(U, W) \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$ form a Markov chain,
4. $U \circlearrowleft (W, Y) \circlearrowleft X$ form a Markov chain.
5. $I(U; Y|Z) = 0$.

Remark 4 Recall that we can use the alternative characterization of $\mathcal{R}^*(\mathbf{D})$, according to Remark 3, where the auxiliary random variables have the Markov structure $U \circlearrowleft V \circlearrowleft W \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$, and increased alphabet sizes. In such a case, conditions 3 and 4 in Theorem 2 are replaced by one condition $U \circlearrowleft W \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$. Furthermore, any Markov chain $U \circlearrowleft W \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$ satisfies conditions 3 and 4 of Theorem 2. We will make use of this fact in the examples in Section 4.

Condition 3 of Theorem 2 implies the Markov structure $U \circlearrowleft Y \circlearrowleft Z$, whereas Condition 5 requires that the roles of Y and Z can be interchanged, i.e., $U \circlearrowleft Z \circlearrowleft Y$. Note that this condition, in general, depends on the distortion measure and on the distortion level (via U), hence it is more general than just requiring that the side information at the two stages be identical: $P_{XY} = P_{XZ}$. We give an example of such source and side information in Section 4.

The random variable V does not play any role in the conditions of Theorem 2. As discussed in the paragraphs following Theorem 1, V is used to encode, at the first stage, information that

can be decoded only at the second stage. It turns out that such a coding scheme is suboptimal for successively refinable sources – if the rate-distortion function can be achieved simultaneously at both distortion levels, then the first stage should transmitt information that can be decoded as a whole by the first decoder.

Proof of Theorem 2. We begin with the necessity part. Assume that (37) holds. From Theorem 1, there exists a triple of random variables (U, V, W) and a pair of maps $f_1 : \mathcal{U} \times \mathcal{Z} \rightarrow \hat{\mathcal{X}}$, $f_2 : \mathcal{W} \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}$, such that

$$(U, V, W) \circlearrowleft X \circlearrowleft Y \circlearrowleft Z \quad (38)$$

is a Markov chain, and moreover

$$R_{X|Z}^*(D_1) \geq I(X; U|Z) + I(X; V|U, Y), \quad (39a)$$

$$\mathbb{E}d(X, f_1(U, Z)) \leq D_1, \quad (39b)$$

and

$$\begin{aligned} R_{X|Y}^*(D_2) &\geq R_{X|Z}^*(D_1) + I(X; W|UVY) \\ &\geq I(X; U|Z) + I(X; VW|UY), \end{aligned} \quad (40a)$$

$$\mathbb{E}d(X, f_2(W, Y)) \leq D_2. \quad (40b)$$

By definition of $R_{X|Z}^*(D_1)$, we conclude that

$$I(X; V|U, Y) = 0 \quad (41)$$

and, in addition, (39a) is satisfied with equality. Thus condition 1 of Theorem 2 holds. The rate bound (40a) can be written as

$$\begin{aligned} R_{X|Y}^*(D_2) &\geq I(X; U|Z) + I(X; VW|UY) \\ &= I(X; U|Y) + I(X; VW|UY) + [I(X; U|Z) - I(X; U|Y)] \\ &\stackrel{(a)}{=} I(X; UVW|Y) + I(U; Y|Z) \\ &\geq I(X; W|Y) + I(U; Y|Z) \\ &\stackrel{(b)}{\geq} R_{X|Y}^*(D_2) + I(U; Y|Z). \end{aligned} \quad (42)$$

where in (a) we used the Markov structure (38), and (b) holds since knowledge of W suffices to satisfy the distortion constraint (see (40b)). Hence, the inequalities in (42) can be replaced by

equalities, and we conclude that

$$I(U; Y|Z) = 0. \quad (43)$$

Therefore, conditions 2 and 5 of the theorem hold. Moreover, from (42) we also conclude that

$$H(X|U, V, W, Y) = H(X|W, Y), \quad (44)$$

which can hold only if $(U, V) \circlearrowleft (W, Y) \circlearrowleft X$. Hence condition 4 is satisfied. Condition 3 follows from (38). This completes the proof of the necessity part.

We proceed to prove the sufficiency part. Assume that there exist a pair of random variables (U, W) and deterministic maps f_1, f_2 , satisfying conditions 1 to 5 of Theorem 2. By the definition of $\mathcal{R}^*(D_1, D_2)$, it remains to show that there exists a random variable V such that (38) holds, and

$$R_{X|Z}^*(D_1) \geq I(X; U|Z) + I(X; V|U, Y), \quad (45)$$

$$R_{X|Y}^*(D_2) - R_{X|Z}^*(D_1) \geq I(X; W|U, V, Y). \quad (46)$$

Indeed, let V be a null random variable (i.e., a constant). Then Condition 3 implies (38). Condition 1 of the theorem implies (45). And for the differential rate we have

$$\begin{aligned} R_{X|Y}^*(D_2) - R_{X|Z}^*(D_1) &\stackrel{(a)}{=} I(X; W|Y) - I(X; U|Z) \\ &\stackrel{(b)}{=} I(X; W|Y) - I(X; U|Y) \\ &= H(X|UY) - H(X|WY) \\ &\stackrel{(c)}{=} H(X|UY) - H(X|UWY) \\ &= I(X; W|UY) \\ &\stackrel{(d)}{=} I(X; W|UVY), \end{aligned} \quad (47)$$

where (a) follows from Conditions 1 and 2 of Theorem 2, (b) follows from Condition 5, (c) follows from Condition 3, and (d) holds since V is independent of all the random variables in the problem. This completes the proof of the sufficiency part. \square

When there is no SI, the conditions of Theorem 2 reduce to those obtained in [7] and [4] for successive refinability without SI. To see this, let Y be deterministic and define $\tilde{f}_1(U) \triangleq f_1(U, Y)$ and $\tilde{f}_2(W) \triangleq f_2(W, Y)$. Conditions 1 and 2 of Theorem 2 become the classical rate-distortion

functions. Conditions 3 and 5 become redundant, and Condition 4 is read $U \circlearrowleft W \circlearrowleft X$. It remains to show that whenever X is successively refinable, \tilde{f}_2 is a one-to-one map (so that $U \circlearrowleft W \circlearrowleft X$ results in $\tilde{f}_1(U) \circlearrowleft \tilde{f}_2(W) \circlearrowleft X$). Notice that for successively refinable source (40b) and (42) imply

$$R(D_2) = I(X; W) \tag{48}$$

$$\mathbb{E}d(X, \tilde{f}_2(W)) \leq D_2, \tag{49}$$

hence \tilde{f}_2 must be one-to-one as otherwise, from the data processing theorem

$$I(X; W) > I(X; \tilde{f}_2(W)), \tag{50}$$

and the coding rate can be reduced to below the rate distortion function $R(D_2)$, which is a contradiction.

It is an easy matter to construct quadruples (U, W, X, Y) where no Markov structure holds for (U, W, X) , yet $U \circlearrowleft (W, Y) \circlearrowleft X$, i.e., Markovity does hold once SI is introduced. Therefore, the presence of Y in condition 4 indicates that a source, that is not successively refinable in the absence of SI, may turn out to be successively refinable when SI is present. We give an example of such a source in Section 4.

4 Examples of Successively Refinable Sources

In this section, we provide three examples of successively refinable sources in the presence of the same SI at both decoders. The first example is the Gaussian source with quadratic distortion measure, the second is the doubly-symmetric binary source with the Hamming distortion measure, and the third example is of a DMS that is not successively refinable in the absence of SI [4], which becomes successively refinable in the presence of SI. In the first two examples, we will use Theorem 2 in conjunction with Remark 4, i.e., we demonstrate a Markov chain $U \circlearrowleft W \circlearrowleft X \circlearrowleft Y$ that meets the conditions of Theorem 2. In addition, for the doubly-symmetric binary source, we construct a degraded SI component, Z , such that $U \circlearrowleft Z \circlearrowleft Y$ holds, therefore it is successively refinable, although the SI at the two stages is not strictly identical.

4.1 The Gaussian Source with Quadratic Distortion

The Wyner–Ziv rate–distortion function of this source is mentioned somewhat briefly in [12]. Before we show that this source is successive refinable, we first describe, in some more detail, the calculation of this rate–distortion function, and then generalize it into two stages. Let $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N \sim \mathcal{N}(0, \sigma_N^2)$ be independent, $Y = X + N$, and let the distortion measure be quadratic, i.e., $d(x, \hat{x}) = (x - \hat{x})^2$. To calculate the Wyner–Ziv rate–distortion function $R_{X|Y}^*(D)$, consider the decomposition² of X as $X = W + S$, where $W \sim \mathcal{N}(0, \sigma_W^2)$ and $S \sim \mathcal{N}(0, \sigma_S^2)$ are independent, thus, $\sigma_X^2 = \sigma_W^2 + \sigma_S^2$. It is this decomposition of σ_X^2 that controls the tradeoff between rate and distortion. Considering W as the auxiliary random variable, one readily obtains, from the Markovity of $W \circlearrowleft X \circlearrowleft Y$, the following expression for the Wyner–Ziv rate:

$$\begin{aligned} I(X; W|Y) &= I(X; W) - I(Y; W) \\ &= \frac{1}{2} \log \left(1 + \frac{\sigma_N^2}{\sigma_S^2} \right) - \frac{1}{2} \log \left(1 + \frac{\sigma_N^2}{\sigma_X^2} \right). \end{aligned} \quad (51)$$

Thus, to minimize $I(X; W|Y)$ over all choices of W , in this class, one should maximize σ_S^2 under the distortion constraint $\min_f \mathbb{E}(X - f(W, Y))^2 \leq D$. Since (W, X, Y) are jointly Gaussian, the best estimator f is linear, i.e., $f(W, Y) = \alpha W + \beta Y$. Upon solving for the optimum coefficients, α and β , the following relation is found between σ_S^2 and D :

$$\sigma_S^2 = \begin{cases} \frac{\sigma_N^2 D}{\sigma_N^2 - D} & D \leq \sigma_{X|Y}^2 \\ \sigma_X^2 & D > \sigma_{X|Y}^2 \end{cases} \quad (52)$$

where

$$\sigma_{X|Y}^2 = \frac{\sigma_X^2 \sigma_N^2}{\sigma_X^2 + \sigma_N^2}. \quad (53)$$

On substituting this expression into (51), one obtains

$$I(X; W|Y) = \begin{cases} \frac{1}{2} \log \frac{\sigma_{X|Y}^2}{D} & D \leq \sigma_{X|Y}^2 \\ 0 & D > \sigma_{X|Y}^2 \end{cases} \quad (54)$$

²In this decomposition, W can be represented as $W = \theta X + S'$, where $\theta = \sigma_W^2 / \sigma_X^2$ and $S' \sim \mathcal{N}(0, \sigma_W^2 (\sigma_X^2 - \sigma_W^2) / \sigma_X^2)$ is independent of X .

whose optimality, as a solution to the Wyner–Ziv problem, is evident from the fact this is also the conditional rate–distortion function when Y is available to the encoder as well.

Turning now to two stages, let W be the auxiliary random variable corresponding to distortion level $D = D_2$, and let us further decompose W as the sum of two independent Gaussian random variables $W = U + T$, U being the other auxiliary random variable of Theorem 2, thus keeping the Markov structure $U \circlearrowleft W \circlearrowleft X \circlearrowleft Y$ (cf. Fig. 2). Since W achieves $R_2 = I(X; W|Y) = R_{X|Y}^*(D_2)$, it remains to show that U achieves $R_1 = I(X; U|Y) = R_{X|Y}^*(D_1)$ for a proper decomposition of W . Since $X = U + S + T$, where U , S , and T are all independent, clearly, the sum $(S + T)$ now plays the previous role of S , and so the appropriate decomposition of W is readily obtained by choosing the variance, σ_T^2 , of T such that

$$\sigma_S^2 + \sigma_T^2 = \begin{cases} \frac{\sigma_N^2 D_1}{\sigma_N^2 - D_1} & D_1 \leq \sigma_{X|Y}^2 \\ \sigma_{X|Y}^2 & D_1 > \sigma_{X|Y}^2, \end{cases} \quad (55)$$

which is always possible since (55) is never smaller than (52) with $D = D_2$. The rate $R_1 = I(X; U|Y)$ will then be as in (54) with $D = D_1$, namely, $R_{X|Y}^*(D_1)$. Thus, we have shown that the Gaussian source with noisy SI is successively refinable with respect to the quadratic distortion measure.

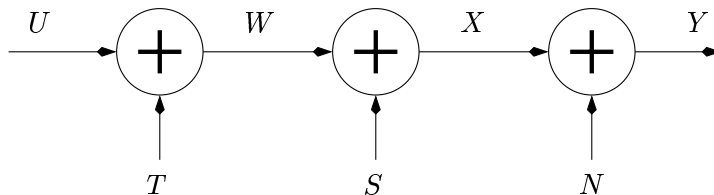


Figure 2: Successive refinement with side information for the Gaussian source with the quadratic distortion measure.

4.2 The Doubly Symmetric Binary Source with Hamming Distortion

The Wyner–Ziv rate–distortion function of the doubly–symmetric binary source with the Hamming distortion measure is calculated in detail in [12]. Roughly speaking, its calculation (as well as its extension to two stages) is analogous to the Gaussian example of subsection 4.1, with the addition operation being replaced by modulo-two addition. It is, however, somewhat more involved because

beyond a certain distortion level, hereafter denoted by D_c , time-sharing with the zero-rate working point must be employed.

As in subsection 4.1, for the sake of completeness, we begin with a description of the system of random variables that corresponds to the solution of the Wyner–Ziv problem of this source, in the ordinary, single-stage setting.

All random variables in this subsection are binary and their alphabet is $\{0, 1\}$. Let X be the binary symmetric source (BSS) and let $Y = X \oplus N$, where \oplus denotes modulo-two addition, and N is an independent random variable with $\Pr\{N = 1\} = p_0$. We are interested in $R_{X|Y}^*(D)$. Define the function

$$g(D) = \begin{cases} h(p_0 * D) - h(D) & 0 \leq D < p_0 \\ 0 & D = p_0 \end{cases} \quad (56)$$

where $h(\cdot)$ is the binary entropy function, and $*$ denotes binary convolution, i.e., $\alpha * \beta = \alpha(1 - \beta) + \beta(1 - \alpha)$. In [12], it is shown that $R_{X|Y}^*(D) = g^*(D)$, the lower convex envelope of $g(D)$. The distortion level D_c is defined as the largest distortion level for which $g^*(D)$ still agrees with $g(D)$, namely, the solution to the equation

$$\frac{g(D_c)}{D_c - p_0} = g'(D_c), \quad (57)$$

where g' is the derivative of g .

We now define the auxiliary random variable as follows. Let S be a random variable, independent of (X, N) , with $\Pr\{S = 1\} = \min\{D, D_c\}$, and define $W_1 = X \oplus S$. Further, let B_1 be independent of (X, N, S) with

$$\Pr\{B = 1\} = \frac{p_0 - \max\{D, D_c\}}{p_0 - D_c} \quad (58)$$

and let $W_2 = B_1 \cdot W_1$. The role of B_1 is to create the time-sharing that is needed for distortion levels above D_c . Finally, define the auxiliary random variable $W = (W_2, B_1)$, and let

$$f(W, Y) = f(W_2, B_1, Y) = B_1 \cdot W_2 + (1 - B_1) \cdot Y. \quad (59)$$

It is then straightforward to show (cf. [12, Section II]) that for the Hamming distortion measure

$$\mathbb{E}d(X, f(W, Y)) = D, \quad (60)$$

and that

$$I(X; W) - I(Y; W) = g^*(D) = R_{X|Y}^*(D). \quad (61)$$

Moving on to two stages of successive refinement, consider the above construction of the auxiliary random variable W with correspondence to distortion level $D = D_2$. We next describe the construction of the additional auxiliary random variable U corresponding to the lower rate, R_1 (cf. Fig. 3). Let T be a random variable, independent of (X, N, S, B_1) such that

$$\Pr\{T = 1\} * \Pr\{S = 1\} = \min\{D_1, D_c\} \quad (62)$$

and let $U_1 = W_2 \oplus T$. Let B_2 be independent of (X, N, S, B_1, T) with

$$\Pr\{B_2 = 1\} = \frac{p_0 - \max\{D_c, D_1\}}{p_0 - \max\{D_c, D_2\}}. \quad (63)$$

Finally, defining $U = (U_2, B_1 \cdot B_2)$, we clearly have $U \ominus W \ominus X \ominus Y$. Letting,

$$f'(U, Y) = B_1 \cdot B_2 \cdot U_2 + (1 - B_1 \cdot B_2)Y. \quad (64)$$

then, similarly as above, it is not difficult to show that

$$\mathbb{E}d(X, f'(U, Y)) = D_1, \quad (65)$$

and that

$$I(X; U) - I(Y; U) = g^*(D_1). \quad (66)$$

Thus, the doubly symmetric binary source is successively refinable with respect to the Hamming distortion measure.

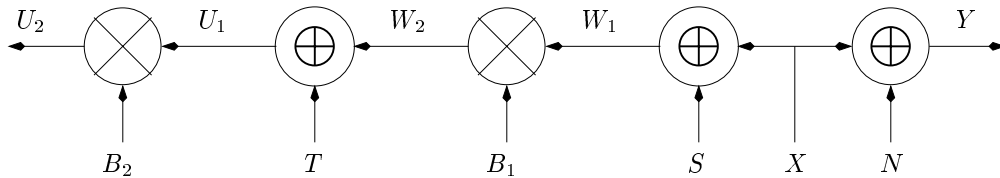


Figure 3: Successive refinement for the doubly symmetric binary source with Hamming distortion measure.

For this source, we now extend the example to include SI which is not strictly identical at the two stages, and yet the source is successively refinable. It is convenient to re-define X to be a BSS

taking values in $\{-1, 1\}$. We define the SI \tilde{Y} to be

$$\tilde{Y} = X \cdot \eta + Q \tag{67}$$

where η is taking values in $\{-1, 1\}$, $P(\eta = -1) = \alpha$, and Q is a real random variable uniformly distributed in the interval $[-1, 1]$. Next, let the degraded random variable Z be defined as

$$Z = \text{sign}(\tilde{Y}). \tag{68}$$

Thus Z is a quantization of \tilde{Y} . Moreover, note that the random variable $X \cdot \eta$ is equivalent to the random variable Y defined above, and that Q does not change its sign. Therefore, although Z is the degraded component and \tilde{Y} and Z are not identical, we also have $U \circlearrowleft X \circlearrowleft Z \circlearrowleft \tilde{Y}$. Hence $I(U; Y|Z) = 0$ and the source is successively refinable with degraded side information.

4.3 The Equitz–Cover Source

Equitz and Cover [4] demonstrate a source that is not successively refinable in the ordinary setting, without SI (see also [3] for more details). The Equitz–Cover example is of a ternary DMS with $\mathcal{X} = \hat{\mathcal{X}} = \{1, 2, 3\}$, a distribution given by $P_X(2) = p$, $P_X(1) = P_X(3) = (1 - p)/2$, for some $p \in [0, 1]$, and a difference distortion measure $d(x, \hat{x}) = \rho(x - \hat{x})$ where ρ is convex and satisfies $\rho(0) < \rho(z)$ whenever $z \neq 0$. It turns out that for some values of D_1 and D_2 , the Markov condition does not hold and thus the source is not successively refinable.

We now show a simple example where in the presence of identical SI, this source becomes successively refinable. Consider a SI variable defined according to

$$Y = \begin{cases} 1 & X \neq 2 \\ 2 & X = 2 \end{cases} \tag{69}$$

Since Y is given by a deterministic function of X , the encoder can reproduce a copy of Y , and so, the Wyner–Ziv rate distortion function must coincide with the conditional rate–distortion function of X given Y at both encoder and decoder. Now, when $Y = 2$, $X \equiv Y$, and so, X is available at the decoder, error–free, without transmission at all. When $Y = 1$, the conditional distribution of X is according to $P_X(1) = P_X(3) = 1/2$, which is actually again an Equitz–Cover source with $p = 0$. Although the reproduction alphabet $\hat{\mathcal{X}}$ is formally still ternary, it is shown in [3, Lemma 4.1.4] that the test channel corresponding to the rate–distortion function of this source induces

$\Pr\{\hat{X} = 2\} = 0$, thus, effectively, the reproduction alphabet reduces to $\hat{\mathcal{X}} = \{1, 3\}$. In other words, given $Y = 1$, we have at hand the rate–distortion problem of the BSS with respect to the Hamming distortion measure, which in turn is successively refinable [4].

5 Multistage Successive Coding

The problem of successive coding can be naturally extended to any finite number of steps. Let $(X, Y_K, Y_{K-1}, \dots, Y_1)$ be a memoryless, stochastically degraded source, i.e.,

$$X \oplus Y_K \oplus Y_{K-1} \oplus \dots \oplus Y_1, \quad (70)$$

where Y_k is taking values in a finite set \mathcal{Y}_k , $1 \leq k \leq K$. Denote by $P_{XY_K Y_{K-1} \dots Y_1}$ its generic distribution, and by P_{XY_k} , $1 \leq k \leq K$, the corresponding marginals. Let $\mathbf{D} = (D_1, D_2, \dots, D_K)$ be a vector of distortion levels. The encoder–decoder pair of the k 'th stage should provide a description of the source with distortion level not exceeding D_k . To this end, the decoder of the k 'th stage, $1 \leq k \leq K$, is provided with the side information vector \mathbf{Y}_k , where

$$\mathbf{Y}_k = (Y_{k,1}, Y_{k,2}, \dots, Y_{k,n}). \quad (71)$$

In addition, each decoder has access to the codewords sent to all its predecessors. We skip the formal definition of a multistage successive code, as it is a straightforward extension of Definition 1 in Section 2. We use $\mathcal{R}_K(\mathbf{D})$ and $\mathcal{R}_{K,i}(\mathbf{D})$ to denote the region of achievable rates for stochastically degraded SI and identical SI, respectively, as in the two stages case.

Let a vector $\mathbf{D} = (D_1, \dots, D_K)$ of distortion levels be given. Define $\mathcal{R}_K^*(\mathbf{D})$ to be the set of all vectors of rates (R_1, R_2, \dots, R_K) for which there exists a collection of $K(K+1)/2$ random variables $\{V_{k,l}, 1 \leq k \leq K, k \leq l \leq K\}$, where $V_{k,l}$ is taking values in a finite set $\mathcal{V}_{k,l}$, such that the following conditions are satisfied

1. $(\{V_{k,l}, 1 \leq k \leq K, k \leq l \leq K\}) \oplus X \oplus Y_K \oplus \dots \oplus Y_1$ is a Markov chain.
2. There exist deterministic maps

$$f_k : \mathcal{V}_{k,k} \times \mathcal{Y}_k \rightarrow \hat{\mathcal{X}}, \quad 1 \leq k \leq K, \quad (72)$$

such that

$$\mathbb{E}d(X, f_k(V_{k,k}, Y_k)) \leq D_k, \quad 1 \leq k \leq K. \quad (73)$$

3. The rates (R_1, R_2, \dots, R_K) satisfy

$$R_1 \geq I(X; V_{1,1}|Y_1) + \sum_{l=2}^K I(X; V_{1,l}|V_{1,1}, V_{1,2}, \dots, V_{1,l-1}, Y_l) \quad (74)$$

$$\begin{aligned} R_k - R_{k-1} &\geq I(X; V_{k,k}|\{V_{i,j}, 1 \leq i < k, 1 \leq j \leq k\}, Y_k) \\ &+ \sum_{l=k+1}^K I(X; V_{k,l}|\{V_{i,j}, 1 \leq i \leq k, i \leq j \leq k\}, Y_k), \quad 2 \leq k \leq K. \end{aligned} \quad (75)$$

Recall that in the case of two stages, the first stage transmits information that can be decoded only with the better SI Y . In the single letter characterization of the region of achievable rates, this information was presented by the random variable V . Similarly, in the multistage case, the encoder of the k 'th stage can transmit information to all the successor decoders. Thus, interpreting the above characterization, for $1 \leq k \leq K$, and $k \leq l \leq K$, the random variable $V_{k,l}$ stands for the information encoded by the k 'th encoder, to be decoded by the l 'th decoder, using Y_l . Since the source is stochastically degraded, $V_{k,l}$ can be decoded also by any of the successor decoders, $l+1, \dots, K$.

Theorem 3 For any discrete, memoryless, stochastically degraded $K+1$ source (X, Y_K, \dots, Y_1) ,

$$\mathcal{R}_K^*(\mathbf{D}) \subseteq \mathcal{R}_K(\mathbf{D}).$$

Proof. The proof employs a hierarchy of random codes, as in the proof of Theorem 1. Although technically involved, it is a straightforward extension of the proof presented in Section 6, and is omitted. \square

Unfortunately, we have not been able to prove a converse of Theorem 3 for general stochastically degraded source. However, we have been able to prove a full coding theorem for identical SI (i.e., $P_{XY_k} = P_{XY} \forall k$), in which case the achievable rate region admits a much simpler form. Define $\mathcal{R}_{K,i}^*(\mathbf{D})$ to be the set of all vectors of rates (R_1, R_2, \dots, R_K) for which there exists a K -tuple of random variables (U_1, U_2, \dots, U_K) , where U_k is taking values in a finite set \mathcal{U}_k , such that the following conditions are satisfied

1. $(U_1, U_2, \dots, U_K) \circlearrowleft X \circlearrowleft Y$ is a Markov chain.

2. There exist deterministic maps

$$f_k : \mathcal{U}_k \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}, \quad 1 \leq k \leq K, \quad (76)$$

such that

$$\mathbb{E}d(X, f_k(U_k, Y)) \leq D_k, \quad 1 \leq k \leq K. \quad (77)$$

3. The rates (R_1, R_2, \dots, R_K) satisfy

$$R_1 \geq I(X; U_1|Y), \quad (78)$$

$$R_k - R_{k-1} \geq I(X; U_k|U_1, U_2, \dots, U_{k-1}, Y_k) \quad \text{for } k = 2, 3, \dots, K. \quad (79)$$

Then we have

Theorem 4 For any discrete, memoryless, joint source (X, Y) ,

$$\mathcal{R}_{K,i}(\mathbf{D}) = \mathcal{R}_{K,i}^*(\mathbf{D}).$$

The proof is presented in Section 6.

A joint source (X, Y) is K -step successively refinable for a given vector of distortion levels $\mathbf{D} = (D_1, D_2, \dots, D_K)$, if

$$(R_{X|Y}^*(D_1), R_{X|Y}^*(D_2), \dots, R_{X|Y}^*(D_K)) \in \mathcal{R}_{K,i}(\mathbf{D}). \quad (80)$$

A by-product of Theorem 4, is a set of necessary and sufficient conditions for (80) to be satisfied.

Theorem 5 A source X with identical side information Y is K -steps successively refinable with distortion levels $\mathbf{D} = (D_1, D_2, \dots, D_K)$, if and only if there exist a K -vector of random variables (U_1, U_2, \dots, U_K) , and K deterministic functions $f_k : \mathcal{U}_k \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}$, $1 \leq k \leq K$, such that the following conditions simultaneously hold:

1. $R_{X|Y}^*(D_k) = I(X; U_k|Y)$ and $\mathbb{E}d(X, f_k(U_k, Y)) \leq D_k$, $1 \leq k \leq K$
2. $(U_1, \dots, U_K) \oplus X \oplus Y$
3. $(U_1, U_2, \dots, U_{k-1}) \oplus (U_k, Y) \oplus X$, $k = 2, 3, \dots, K$.

It is worth pointing out that, in the examples given in Secs. 4.1 and 4.2 (Gaussian source with quadratic distortion measure, and doubly symmetric binary source with Hamming distortion measure, resp.), the construction of auxiliary random variables U and W can be extended to any number of stages, thus these sources are K -step successively refinable.

Proof of Theorem 5. We start with necessity. Assume that (80) holds. From Theorem 4, there exists a K -tuple (U_1, U_2, \dots, U_K) and K deterministic maps $f_k : \mathcal{U}_k \times \mathcal{Y} \rightarrow \hat{\mathcal{X}}$, $1 \leq k \leq K$, such that

$$R_{X|Y}^*(D_1) \geq I(X; U_1|Y), \quad (81a)$$

$$\mathbb{E}d(X, f_1(U_1, Y)) \leq D_1, \quad (81b)$$

similarly,

$$R_{X|Y}^*(D_2) \geq R_{X|Y}^*(D_1) + I(X; U_2|U_1, Y) \geq I(X; U_1 U_2|Y), \quad (82a)$$

$$\mathbb{E}d(X, f_2(U_2, Y)) \leq D_2, \quad (82b)$$

and, in general

$$R_{X|Y}^*(D_k) \geq I(X; U_1 \dots U_k|Y), \quad (83a)$$

$$\mathbb{E}d(X, f_k(U_k, Y)) \leq D_k, \quad 1 \leq k \leq K. \quad (83b)$$

By the definition of $R_{X|Y}^*(D)$, (81a) must be satisfied with equality. For (83) we have

$$R_{X|Y}^*(D_k) \geq I(X; U_1 \dots U_k|Y) \geq I(X; U_k|Y) \geq R_{X|Y}^*(D_k), \quad (84)$$

where the last inequality is due to the fact that knowledge of U_k and Y suffices to satisfy the distortion constraint (see (83)). Hence equality in (84) must hold, and Conditions 1 and 2 of the theorem hold. For equality hold in (84), we must have

$$H(X|U_1 U_2 \dots U_{k-1} U_k, Y) = H(X|U_k, Y), \quad k = 2, 3, \dots, K, \quad (85)$$

implying Condition 3 of the theorem. This concludes the proof of necessity.

To prove sufficiency, assume that there exist a K -tuple of random variables (U_1, \dots, U_K) , and K independent maps f_1, \dots, f_K such that Conditions 1 to 3 of the theorem are met. By the definition of $\mathcal{R}_{K,i}^*(\mathbf{D})$, it remains to show that the differential rates satisfy the bounds

$$R_{X|Y}^*(D_k) - R_{X|Y}^*(D_{k-1}) \geq I(X; U_k|U_1 U_2 \dots U_{k-1}, Y), \quad k = 2, 3, \dots, K. \quad (86)$$

Indeed, we can write

$$\begin{aligned}
R_{X|Y}^*(D_k) - R_{X|Y}^*(D_{k-1}) &= I(X; U_k|Y) - I(X; U_{k-1}|Y) \\
&= H(X|U_{k-1}, Y) - H(X|U_k, Y) \\
&\stackrel{(a)}{=} H(X|U_{k-1}, Y) - H(X|U_1U_2 \dots U_k, Y) \\
&\stackrel{(b)}{=} H(X|U_1U_2 \dots U_{k-1}, Y) - H(X|U_1U_2 \dots U_k, Y) \\
&= I(X; U_k|U_1U_2 \dots U_{k-1}, Y), \tag{87}
\end{aligned}$$

where (a) and (b) are due to Condition 3 of the theorem. This concludes to proof of sufficiency. \square

6 Proofs of Theorems 1 and 4

6.1 Proof of Theorem 1

Converse part: Assume that we have an (n, M_1, M_2, D_1, D_2) SR code for the source X with SI (Y, Z) , as in Definition 1. We will show the existence of a triple (U, V, W) that satisfies conditions 1-4 in the definition of $\mathcal{R}^*(\mathbf{D})$. Denote $T_i = \phi_i(X^n)$, $i = 1, 2$. Then

$$\begin{aligned}
nR_1 &\geq H(T_1) \geq I(X^n; T_1|Z^n) = I(X^n; T_1Y^n|Z^n) - I(X^n; Y^n|T_1Z^n) \\
&= \sum_{i=1}^n \left[I(X_i; T_1Y^n|X^{i-1}Z^n) - I(X^n; Y_i|T_1Z^nY^{i-1}) \right]. \tag{88}
\end{aligned}$$

For notational convenience, we denote $Z^{i-1}Z_{i+1}^n = Z^{n \setminus i}$, and use a similar notation for X and Y . Since X_iZ_i and $X^{i-1}Z^{n \setminus i}$ are independent, we have, for the first term in the summand of (88):

$$\begin{aligned}
I(X_i; T_1Y^n|X^{i-1}Z^n) &= H(X_i|Z_iX^{i-1}Z^{n \setminus i}) - H(X_i|Z_iX^{i-1}Z^{n \setminus i}T_1Y^n) \\
&= H(X_i|Z_i) - H(X_i|Z_iX^{i-1}Z^{n \setminus i}T_1Y^n) \\
&= I(X_i; X^{i-1}Z^{n \setminus i}T_1Y^n|Z_i). \tag{89}
\end{aligned}$$

Next, due to the Markov structure

$$Y_i \ominus (X_iZ_i) \ominus (X^{n \setminus i}T_1Y^{i-1}Z^{n \setminus i}) \tag{90}$$

we have, for the second term in the summand of (88):

$$\begin{aligned}
I(X^n; Y_i | T_1 Z^n Y^{i-1}) &= H(Y_i | T_1 Z^n Y^{i-1}) - H(Y_i | X^n T_1 Z^n Y^{i-1}) \\
&= H(Y_i | T_1 Z^n Y^{i-1}) - H(Y_i | X_i T_1 Z^n Y^{i-1}) \\
&= I(X_i; Y_i | T_1 Z^n Y^{i-1}).
\end{aligned} \tag{91}$$

Substituting (89) and (91) in (88), we obtain

$$\begin{aligned}
nR_1 &\geq \sum_{i=1}^n \left[I(X_i; X^{i-1} Z^{n \setminus i} T_1 Y^n | Z_i) - I(X_i; Y_i | T_1 Z^n Y^{i-1}) \right] \\
&= \sum_{i=1}^n \left[I(X_i; Z^{n \setminus i} T_1 Y^{i-1} | Z_i) + I(X_i; X^{i-1} Y_i^n | Z_i T_1 Z^{n \setminus i} Y^{i-1}) - I(X_i; Y_i | T_1 Z^n Y^{i-1}) \right] \\
&= \sum_{i=1}^n \left[I(X_i; T_1 Z^{n \setminus i} Y^{i-1} | Z_i) + I(X_i; X^{i-1} Y_{i+1}^n | Z_i Y_i T_1 Z^{n \setminus i} Y^{i-1}) \right].
\end{aligned} \tag{92}$$

The Markovity of $X \ominus Y \ominus Z$ implies

$$Z_i \ominus Y_i \ominus (X_i T_1 Z^{n \setminus i} Y^{i-1}), \tag{93}$$

and we have for the second term in (92)

$$\begin{aligned}
&I(X_i; X^{i-1} Y_{i+1}^n | T_1 Z^n Y^i) \\
&= H(X_i | T_1 Z^n Y^i) - H(X_i | T_1 Z^n Y^n X^{i-1}) \\
&= H(X_i Z_i | T_1 Z^{n \setminus i} Y^i) - H(Z_i | T_1 Z^{n \setminus i} Y^i) - H(X_i | T_1 Z^n Y^n X^{i-1}) \\
&= H(Z_i | X_i T_1 Z^{n \setminus i} Y^i) + H(X_i | T_1 Z^{n \setminus i} Y^i) - H(Z_i | T_1 Z^{n \setminus i} Y^i) - H(X_i | T_1 Z^n Y^n X^{i-1}) \\
&\stackrel{(a)}{=} H(X_i | T_1 Z^{n \setminus i} Y^i) - H(X_i | T_1 Z^n Y^n X^{i-1}) = I(X_i; Z_i Y_{i+1}^n X^{i-1} | T_1 Z^{n \setminus i} Y^i) \\
&\geq I(X_i; Y_{i+1}^n X^{i-1} | T_1 Z^{n \setminus i} Y^i)
\end{aligned} \tag{94}$$

where (93) was used in (a). Substituting (94) in (92), we get

$$nR_1 \geq \sum_{i=1}^n \left[I(X_i; T_1 Z^{n \setminus i} Y^{i-1} | Z_i) + I(X_i; Y_{i+1}^n X^{i-1} | T_1 Z^{n \setminus i} Y^i) \right]. \tag{95}$$

Before defining the auxiliary random variables, we bound the refinement rate from below as follows:

$$\begin{aligned}
n(R_2 - R_1) &\geq H(T_2 | T_1) \geq H(T_2 | T_1 Z^n Y^n) \geq I(X^n; T_2 | T_1 Z^n Y^n) \\
&= \sum_{i=1}^n I(X_i; T_2 | X^{i-1} T_1 Z^n Y^n).
\end{aligned} \tag{96}$$

Define the random variables

$$U_i = T_1 Z^{n \setminus i} Y^{i-1} \quad (97)$$

$$V_i = X^{i-1} Y_{i+1}^n U_i \quad (98)$$

$$W_i = T_2 V_i. \quad (99)$$

With these definitions, we have the Markov structure

$$U_i \ominus V_i \ominus W_i \ominus X_i \ominus Y_i \ominus Z_i \quad (100)$$

and the bounds (95) and (96) become

$$R_1 \geq \frac{1}{n} \sum_{i=1}^n [I(X_i; U_i | Z_i) + I(X_i; V_i | U_i, Y_i)] \quad (101)$$

$$R_2 - R_1 \geq \frac{1}{n} \sum_{i=1}^n I(X_i; W_i | U_i, V_i, Y_i) \quad (102)$$

where we have used (100) to drop the conditioning on Z_i in (102).

Let J be a random variable, independent of X , Y , and Z , and uniformly distributed over the set $\{1, 2, \dots, n\}$. Define the random variables $U = (J, U_J)$, $V = (J, V_J)$, and $W = (J, W_J)$. The Markov relations (100) still hold, that is

$$U \ominus V \ominus W \ominus X \ominus Y \ominus Z, \quad (103)$$

and therefore the condition 1 in the definition of $\mathcal{R}^*(\mathbf{D})$ is satisfied. We proceed to show the existence of functions f_1, f_2 satisfying (12), (13). Denote by $\psi_{i,k}$ the output of the i -th decoder at time k , $i = 1, 2$, $1 \leq k \leq n$. The random variable U contains $\phi_1(X^n) Z^{n \setminus J}$. Similarly, W contains $\phi_1(X^n) \phi_2(X^n) Z^{n \setminus J} Y^{n \setminus J}$. Therefore, let us choose the functions f_1 and f_2 as follows

$$f_1(U, Z) = \psi_{1,J}(\phi_1(X^n), Z^n) \quad (104)$$

$$f_2(W, Y) = \psi_{2,J}(\phi_1(X^n), \phi_2(X^n), Y^n). \quad (105)$$

Then, for the distortions we have

$$\mathbb{E}d(X, f_1(U, Z)) = \frac{1}{n} \sum_{j=1}^n \mathbb{E}d(X, \psi_{1,j}(\phi_1(X^n), Z^n)) \leq D_1 \quad (106)$$

$$\mathbb{E}d(X, f_2(W, Y)) = \frac{1}{n} \sum_{j=1}^n \mathbb{E}d(X, \psi_{2,j}(\phi_1(X^n), \phi_2(X^n), Y^n)) \leq D_2 \quad (107)$$

Hence, condition 2 in the definition of $\mathcal{R}^*(\mathbf{D})$ is satisfied. To prove that condition 4 of that definition holds, we have to show that the bounds (101) and (102) can be written in a single letter form with U , V , and W . The following chain of equalities holds

$$\begin{aligned}
I(X; U|Z) &= H(U|Z) - H(U|XZ) = H(U|Z) - H(U|X) \\
&= H(U) - H(U|X) - (H(U) - H(U|Z)) \\
&= I(U; X) - I(U; Z) \\
&= H(X) - H(X|U) - H(Z) + H(Z|U) \\
&= H(X) - H(X|J, U_J) - H(Z) + H(Z|J, U_J) \\
&= \frac{1}{n} \sum_{i=1}^n H(X_i) - \frac{1}{n} \sum_{i=1}^n H(X_i|U_i) - \frac{1}{n} \sum_{i=1}^n H(Z_i) + \frac{1}{n} \sum_{i=1}^n H(Z_i|U_i) \\
&= \frac{1}{n} \sum_{i=1}^n I(X_i; U_i|Z_i)
\end{aligned} \tag{108}$$

where the last equality is due to (100). In a similar manner, we get

$$\begin{aligned}
I(X; V|U, Y) &= I(X; JV_J|JU_J, Y) = I(X; V_J|JU_J, Y) \\
&= H(X|JU_J, Y) - H(X|JU_JV_J, Y) \\
&= \frac{1}{n} \sum_{i=1}^n H(X_i|i, U_i, Y_i) - \frac{1}{n} \sum_{i=1}^n H(X_i|i, U_i, V_i, Y_i) \\
&= \frac{1}{n} \sum_{i=1}^n I(X_i; V_i|U_i, Y_i).
\end{aligned} \tag{109}$$

In view of (108), (109), the bound (101) can be written as

$$R_1 \geq I(X; U|Z) + I(X; V|U, Y). \tag{110}$$

In a similar manner, it is shown that (102) can be written as

$$R_2 - R_2 \geq I(X; W|U, V, Y). \tag{111}$$

To complete the proof of the converse part, it remains to prove that condition 3 in the definition of $\mathcal{R}^*(\mathbf{D})$ holds, namely, to show that the alphabets of the random variables U , V , and W , can be limited, without loss of generality, as in eqs. (14), (15), and (16), respectively. To this end, we invoke the support lemma (cf. [2, p. 310]), which results from the Carathéodory theorem.

According to this lemma, given k real-valued, continuous functionals f_j , $j = 1, \dots, k$, on the set $\mathcal{P}(\mathcal{X})$ of probability distributions over an alphabet \mathcal{X} , and given any probability measure μ on the Borel σ -algebra of $\mathcal{P}(\mathcal{X})$, there exist k elements Q_1, \dots, Q_k of $\mathcal{P}(\mathcal{X})$ and k non-negative reals, $\alpha_1, \dots, \alpha_k$, that sum to unity, such that for every $j = 1, \dots, k$:

$$\int_{\mathcal{P}(\mathcal{X})} f_j(Q) \mu(dQ) = \sum_{i=1}^k \alpha_i f_j(Q_i). \quad (112)$$

Before we actually apply the support lemma, we first rewrite the relevant conditional mutual informations and the distortion functions (that appear in the definition of $\mathcal{R}^*(\mathbf{D})$) in a more convenient form for the use of this lemma, by taking advantage of the Markov structure. As for the first term, $I(X; U|Z)$, in the lower bound to R_1 , we have:

$$\begin{aligned} I(X; U|Z) &= H(U|Z) - H(U|X, Z) \\ &= H(U|Z) - H(U|X) \\ &= H(U) + H(Z|U) - H(Z) - H(U) - H(X|U) + H(X) \\ &= H(X) - H(Z) + H(Z|U) - H(X|U). \end{aligned} \quad (113)$$

For the second term in the lower bound to R_1 , we have

$$\begin{aligned} I(X; V|U, Y) &= H(X|U, Y) - H(X|U, V, Y) \\ &= \mathbb{E} \log P(X|U, V, Y) - \mathbb{E} \log P(X|U, Y) \\ &= \mathbb{E} \log \left[\frac{P(U, V, X, Y)}{P(U, V, Y)} \right] - \mathbb{E} \log \left[\frac{P(U, X, Y)}{P(U, Y)} \right] \\ &= \mathbb{E} \log \left[\frac{P(U, V)P(X|U, V)P(Y|X)}{P(U, V)P(Y|U, V)} \right] - \mathbb{E} \log \left[\frac{P(U)P(X|U)P(Y|X)}{P(U)P(Y|U)} \right] \\ &= H(X|U) - H(Y|U) + H(Y|U, V) - H(X|U, V). \end{aligned} \quad (114)$$

Thus, the sum of these two terms is given by:

$$\begin{aligned} I(X; U|Y) + I(X; V|U, Y) &= [H(X) - H(Z)] + [H(Z|U) - H(Y|U)] + \\ &\quad [H(Y|U, V) - H(X|U, V)] \\ &\stackrel{\Delta}{=} \alpha + \beta + \gamma. \end{aligned} \quad (115)$$

In a similar manner, for the lower bound to $R_2 - R_1$, we obtain

$$I(X; W|U, V, Y) = [H(Y|U, V, W) - H(X|U, V, W)] - \gamma \stackrel{\Delta}{=} \delta - \gamma. \quad (116)$$

In order to preserve prescribed values of $I(X; U|Z) + I(X; V|U, Y)$ and $I(X; W|U, V, Y)$, it then suffices to preserve the associated values of $(\beta + \gamma)$ and $(\delta - \gamma)$, since α is a constant that depends only on the given statistics of the source and the SI.

We first invoke the support lemma in order to reduce the alphabet size of U , while preserving the values of $\beta + \gamma$ and $\delta - \gamma$ as well as the first-stage distortion. The alphabets of V and W are still kept intact at this step. Define the following functionals of a generic distribution Q over $\mathcal{X} \times \mathcal{V} \times \mathcal{W}$, where \mathcal{X} is assumed, without loss of generality, to be $\{1, 2, \dots, m\}$, $m \triangleq |\mathcal{X}|$.

$$f_x(Q) = \sum_{v,w} Q(x, v, w), \quad x = 1, \dots, m-1 \quad (117)$$

$$f_m(Q) = A(Q) + B(Q), \quad (118)$$

where

$$A(Q) = \sum_{x,v,w} Q(x, v, w) \sum_{y,z} P(y, z|x) \log \frac{\sum_{x',v',w'} Q(x', v', w') P(y|x')}{\sum_{x',v',w'} Q(x', v', w') P(z|x')}, \quad (119)$$

and

$$B(Q) = \sum_{x,v,w} Q(x, v, w) \sum_y P(y|x) \log \frac{Q(x|v)}{P_Q(y|v)} \quad (120)$$

with

$$Q(x|v) \triangleq \frac{\sum_{w'} Q(x, v, w')}{\sum_{x',w'} Q(x', v, w')} \quad (121)$$

and

$$P_Q(y|v) \triangleq \frac{\sum_{x',w'} P(y|x') Q(x', v, w')}{\sum_{x',w',y'} P(y'|x') Q(x', v, w')}. \quad (122)$$

Next, define

$$f_{m+1}(Q) = \sum_{x,v,w} Q(x, v, w) \sum_y P(y|x) \log \frac{Q(x|v, w)}{P_Q(y|v, w)} - B(Q) \quad (123)$$

with similar definitions (but without the summations over w') for the numerator and the denominator of the argument of the logarithm. Finally, for the distortion constraint, let

$$f_{m+2}(Q) = \sum_z \min_{\hat{x}} \sum_{x,v,w} Q(x, v, w) P(z|x) d(x, \hat{x}). \quad (124)$$

Applying now the support lemma, we find that there exists a random variable U (jointly distributed with X) whose alphabet size is $k = m + 2 = |\mathcal{X}| + 2$ that satisfies simultaneously

$$\sum_u P(U = u) f_x(P(\cdot|u)) = P(x), \quad x = 1, \dots, m - 1 \quad (125)$$

$$\sum_u P(U = u) f_m(P(\cdot|u)) = \beta + \gamma \quad (126)$$

$$\sum_u P(U = u) f_{m+1}(P(\cdot|u)) = \delta - \gamma \quad (127)$$

$$\sum_u P(U = u) f_{m+2}(P(\cdot|u)) = \min_f \mathbb{E}d(X, f(U, Z)). \quad (128)$$

Having found such a random variable U , we now proceed to reduce the alphabet of V in a similar manner, where this time, we have $|\mathcal{X}| \cdot |\mathcal{U}| - 1$ constraints to preserve the joint distribution of (X, U) just defined, and two more constraints for preserving γ and δ (note that V is not involved in β and in any distortion constraint). Thus, the necessary alphabet size of V is upper bounded by

$$|\mathcal{V}| \leq |\mathcal{X}| \cdot |\mathcal{U}| + 1 \leq |\mathcal{X}| \cdot (|\mathcal{X}| + 2) + 1 = (|\mathcal{X}| + 1)^2. \quad (129)$$

Finally, W must preserve the joint distribution of (X, U, V) , plus the value of δ and the second expected distortion $\min_f \mathbb{E}d(X, f(W, Y))$, which means that the needed alphabet size of W does not have to exceed

$$|\mathcal{X}| \cdot |\mathcal{U}| \cdot |\mathcal{V}| + 1 \leq |\mathcal{X}| \cdot (|\mathcal{X}| + 2) \cdot (|\mathcal{X}| + 1)^2 + 1. \quad (130)$$

This completes the proof of the converse part.

Direct Part:

We begin by setting up some notation and mentioning a few basic facts that will be needed hereafter. Given a distribution P_{XU} , we denote by \mathcal{T}_X^δ the set of all n -tuples \mathbf{x} that are δ -typical according to P_X , i.e.,

$$\mathcal{T}_X^\delta = \left\{ \mathbf{x} \in \mathcal{X}^n : \left| \frac{1}{n} N(x|\mathbf{x}) - P_X(x) \right| \leq \delta \quad \forall x \in \mathcal{X}, \right. \\ \left. \text{and } N(x|\mathbf{x}) = 0 \text{ whenever } P_X(x) = 0 \right\}$$

where $N(x|\mathbf{x})$ is the number of occurrences of the letter x in the n -tuple \mathbf{x} . In the sequel, we will use the following well known results [2]. For any $\mathbf{x} \in \mathcal{T}_X^\delta$ and any $\delta' > \delta$

$$\exp[-n(I(X; U) + \epsilon_u)] \leq \sum_{\mathbf{u}: (\mathbf{x}, \mathbf{u}) \in \mathcal{T}_{XU}^{\delta'}} P_U(\mathbf{u}) \leq \exp[-n(I(X; U) - \epsilon)] \quad (131)$$

where $\epsilon = \epsilon(\delta, \delta')$ and $\epsilon_u = \epsilon_u(\delta, \delta')$ both vanish as $\delta, \delta' \rightarrow 0$. Similarly, for any pair $(\mathbf{x}, \mathbf{u}) \in \mathcal{T}_{XU}^{\delta'}$ and any $\delta'' > \delta'$

$$\exp[-n(I(X; V|U) + \eta_{v|u})] \leq \sum_{\mathbf{v}: (\mathbf{x}, \mathbf{u}, \mathbf{v}) \in \mathcal{T}_{XUV}^{\delta''}} P_{V|U}(\mathbf{v}|\mathbf{u}) \leq \exp[-n(I(X; V|U) - \eta)] \quad (132)$$

where $\eta = \eta(\delta', \delta'')$ and $\eta_{v|u} = \eta_{v|u}(\delta', \delta'')$ both vanish as $\delta', \delta'' \rightarrow 0$.

Let a distortion pair $\mathbf{D} = (D_1, D_2)$ be given, and let (U, V, W) satisfy the conditions that define $\mathcal{R}^*(\mathbf{D})$. Recall that this guarantees the existence of functions f_1, f_2 , satisfying (12), (13). Fix an arbitrary $\gamma > 0$ and consider the construction of an hierarchy of codebooks described in the following steps.

Codebook Generation

1. Randomly generate $M_1 = \exp[n(I(X; U) + \gamma)]$ independent codewords $\{\mathbf{u}_i\}$, each of length n , according to $P_U(\cdot)$. Denote

$$\mathcal{A} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{M_1}\}, \quad \mathbf{u}_j \in \mathcal{U}^n, \quad 1 \leq j \leq M_1. \quad (133)$$

2. Let $L_1 = \exp[n(I(Z; U) - 4\gamma)]$, and

$$\begin{aligned} N_1 &= \exp[n(I(X; U) - I(Z; U) + 5\gamma)] \\ &\stackrel{(a)}{=} \exp[n(I(X; U|Z) + 5\gamma)] \\ &= M_1 \exp[-n(I(Z; U) - 4\gamma)] = \frac{M_1}{L_1}, \end{aligned} \quad (134)$$

where (a) above follows from the Markov structure $(U, V, W) \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$. Partition the codebook \mathcal{A} into N_1 bins, each containing L_1 members of \mathcal{A} . Let $\mathcal{A}_B(i)$ denote the elements $\mathbf{u} \in \mathcal{A}$ assigned to bin i , $1 \leq i \leq N_1$.

3. Set $M_2 = \exp[n(I(X; V|U) + \gamma)]$. For any $\mathbf{u} \in \mathcal{A}$, generate a codebook, of length n and size M_2 , according to $P_{V|U}(\mathbf{v}|\mathbf{u})$, where

$$P_{V|U}(\mathbf{v}|\mathbf{u}) = \prod_{l=1}^n P_{V|U}(v_l|u_l). \quad (135)$$

Denote this codebook by $\mathcal{B}(\mathbf{u})$.

4. Let $L_2 = \exp [n(I(Y; V|U) - 4\gamma)]$, and

$$\begin{aligned}
N_2 &= \exp [n(I(X; V|U) - I(Y; V|U) + 5\gamma)] \\
&= \exp [n(I(X; V|U, Y) + 5\gamma)] \\
&= M_2 \exp [-n(I(Y; V|U) - 4\gamma)] = \frac{M_2}{L_2}.
\end{aligned} \tag{136}$$

For each element $\mathbf{u} \in \mathcal{A}$, partition the codebook $\mathcal{B}(\mathbf{u})$ into N_2 bins, each containing L_2 members of $\mathcal{B}(\mathbf{u})$. Let $\mathcal{B}_B(j, \mathbf{u})$ denote the set of elements $\mathbf{v} \in \mathcal{B}(\mathbf{u})$ assigned to bin j , $1 \leq j \leq N_2$.

5. Set $M_3 = \exp [n(I(X; W|U, V) + \gamma)]$. For any pair (\mathbf{u}, \mathbf{v}) such that $\mathbf{u} \in \mathcal{A}$ and $\mathbf{v} \in \mathcal{B}(\mathbf{u})$, generate a codebook of length n and size M_3 , according to $P_{W|UV}$

$$P_{W|UV}(\mathbf{w}|\mathbf{u}, \mathbf{v}) = \prod_{l=1}^n P_{W|UV}(w_l|u_l, v_l). \tag{137}$$

Denote this codebook by $\mathcal{C}(\mathbf{u}, \mathbf{v})$.

6. Similarly to the partitions of \mathcal{A} and $\mathcal{B}(\mathbf{u})$ above, we now partition $\mathcal{C}(\mathbf{u}, \mathbf{v})$. Thus, define $L_3 = \exp [n(I(Y; W|V, U) - 4\gamma)]$, and

$$\begin{aligned}
N_3 &= \exp [n(I(X; W|U, V) - I(Y; W|U, V) + 5\gamma)] \\
&= \exp [n(I(X; W|U, V, Y) + 5\gamma)] \\
&= M_3 \exp [-n(I(Y; W|U, V) - 4\gamma)] = \frac{M_3}{L_3}.
\end{aligned} \tag{138}$$

For each pair (\mathbf{u}, \mathbf{v}) such that $\mathbf{u} \in \mathcal{A}$ and $\mathbf{v} \in \mathcal{B}(\mathbf{u})$, partition the codebook $\mathcal{C}(\mathbf{u}, \mathbf{v})$ into N_3 bins, each containing L_3 members. Let $\mathcal{C}_B(k, \mathbf{u}, \mathbf{v})$ denote the elements $\mathbf{w} \in \mathcal{C}(\mathbf{u}, \mathbf{v})$ assigned to bin k , $1 \leq k \leq N_3$.

Reveal the codebooks and the partitions to the encoder and decoders.

Encoding

Given a source vector \mathbf{x} , the encoding process proceeds along the following steps.

1. First encoder: The encoding map $\phi_1(\mathbf{x})$ consists of two maps,

$$\phi_1(\mathbf{x}) = \phi_{1,1}(\mathbf{x})\phi_{1,2}(\mathbf{x}),$$

i.e., it is a double index, defined as follows. The encoder seeks a vector $\mathbf{u} \in \mathcal{A}$ such that $(\mathbf{x}, \mathbf{u}) \in \mathcal{T}_{XU}^{2\delta}$. If more than one such vector exists in the codebook \mathcal{A} , the first one is chosen. If such a vector does not exist in the codebook \mathcal{A} , a default vector is chosen, say \mathbf{u}_1 , and an error is declared. Denote this vector by $\mathbf{u}(\mathbf{x})$. The second index is set equal to the bin number to which $\mathbf{u}(\mathbf{x})$ belongs, that is

$$\phi_{1,2}(\mathbf{x}) = i \quad \text{if} \quad \mathbf{u}(\mathbf{x}) \in \mathcal{A}_B(i). \quad (139)$$

The mapping $\phi_{1,2}(\cdot)$ thus takes values in the set $\{1, 2, \dots, N_1\}$.

We proceed to the first map, $\phi_{1,1}$. Given \mathbf{x} and a codeword $\mathbf{u}(\mathbf{x}) \in \mathcal{A}$, the encoder seeks a vector $\mathbf{v} \in \mathcal{B}(\mathbf{u}(\mathbf{x}))$ such that $(\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}) \in \mathcal{T}_{XUV}^{3\delta}$. If such vector does not exist, a default vector is chosen, say the first vector in $\mathcal{B}(\mathbf{u}(\mathbf{x}))$, and an error is declared. If more than one exists, the first one is chosen. Denote this vector by $\mathbf{v}(\mathbf{x})$. The map $\phi_{1,1}(\mathbf{x})$ is defined as the bin number to which $\mathbf{v}(\mathbf{x})$ belongs, that is

$$\phi_{1,1}(\mathbf{x}) = j \quad \text{if} \quad \mathbf{v}(\mathbf{x}) \in \mathcal{B}_B(j, \mathbf{u}(\mathbf{x})). \quad (140)$$

Observe that the encoding map ϕ_1 is taking values in a discrete set of size $N_1 \cdot N_2$, where

$$N_1 \cdot N_2 = \exp[n(I(X; V|Z) + I(X; V|U, Y) + 10\gamma)].$$

2. Second encoder (refinement encoder): The second encoder knows \mathbf{x} , $\mathbf{u}(\mathbf{x})$, and $\mathbf{v}(\mathbf{x})$. It seeks a codeword $\mathbf{w} \in \mathcal{C}(\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))$ such that $(\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}), \mathbf{w}) \in \mathcal{T}_{XUVW}^{4\delta}$. If such vector does not exist, a default \mathbf{w} is chosen, say, the first vector in $\mathcal{C}(\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))$, and an error is declared. If more than one exists, the first in the list is chosen. Denote this vector by $\mathbf{w}(\mathbf{x})$. The value assigned to $\phi_2(\mathbf{x})$ is the bin number to which $\mathbf{w}(\mathbf{x})$ belongs, i.e.,

$$\phi_2(\mathbf{x}) = k \quad \text{if} \quad \mathbf{w}(\mathbf{x}) \in \mathcal{C}_B(k, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x})). \quad (141)$$

Observe that the mapping ϕ_2 is taking values in $\{1, 2, \dots, N_3\}$.

The first encoder sends the indices of the two bins, $\phi_{1,1}(\mathbf{x})\phi_{1,2}(\mathbf{x})$, using

$$\log(N_1 \cdot N_2) = n[I(X; U|Z) + I(X; V|U, Y) + 10\gamma] \quad \text{bits}. \quad (142)$$

The refinement encoder sends the index of the third bin using

$$\log N_3 = n[I(X; W|U, V, Y) + 5\gamma] \quad \text{bits}. \quad (143)$$

Decoding

1. The decoder of the first stage has the vector \mathbf{z} as SI. It has access to the double index $\phi_2(\mathbf{x}) = (i, j)$, $1 \leq i \leq N_1$, $1 \leq j \leq N_2$, but in the reconstruction process it utilizes only the first index i . It decodes the vector \mathbf{u} precisely as in the classical Wyner–Ziv decoding procedure [1], [12]. Specifically, the first decoder seeks a unique vector $\mathbf{u} \in \mathcal{A}_B(i)$ such that $(\mathbf{u}, \mathbf{z}) \in \mathcal{T}_{UZ}^{3\delta|\mathcal{X}|}$. Denote this vector by $\hat{\mathbf{u}}(\mathbf{z})$. If there is no vector $\mathbf{u} \in \mathcal{A}_B(i)$ jointly typical with \mathbf{z} , or there is more than one, an arbitrary $\hat{\mathbf{u}}$ is chosen, and an error is declared. The reconstruction vector of the first stage, $\hat{\mathbf{x}}_1$, is given by

$$\hat{\mathbf{x}}_1 = (\hat{x}_{1,1}, \hat{x}_{1,2}, \dots, \hat{x}_{1,n}) \quad (144)$$

where

$$\hat{x}_{1,l} = f_1(\hat{u}_l(\mathbf{z}), z_l). \quad (145)$$

2. The decoder of the second stage has the vector \mathbf{y} as SI. It receives the double-index (i, j) as the first decoder, but utilizes both indices, i and j . In addition, it gets the output of the second encoder, $\phi_2(\mathbf{x}) = k$, $1 \leq k \leq N_3$, and utilizes it too. Thus, it can be described in four stages.

- (a) Decoding \mathbf{u} : Due to the Markov structure $(U, V, W) \circlearrowleft X \circlearrowleft Y \circlearrowleft Z$, the second decoder can do everything the first decoder can. Thus, it looks for a unique vector $\mathbf{u} \in \mathcal{A}_B(i)$ such that

$$(\mathbf{u}, \mathbf{y}) \in \mathcal{T}_{UY}^{3\delta|\mathcal{X}|}. \quad (146)$$

If such vector \mathbf{u} does not exist, or if it is not unique, an arbitrary element in $\mathcal{A}_B(i)$ is chosen. Denote the output of this stage by $\hat{\mathbf{u}}(\mathbf{y})$.

- (b) Decoding \mathbf{v} : In this stage, the second decoder looks for a unique vector $\mathbf{v} \in \mathcal{B}_B(j, \hat{\mathbf{u}}(\mathbf{y}))$ such that

$$(\hat{\mathbf{u}}(\mathbf{y}), \mathbf{v}, \mathbf{y}) \in \mathcal{T}_{UVY}^{4\delta|\mathcal{X}|}. \quad (147)$$

Denote this vector by $\hat{\mathbf{v}}(\mathbf{y})$. As usual, if such a vector does not exist, or there is more than one such vector in $\mathcal{B}_B(j, \hat{\mathbf{u}}(\mathbf{y}))$, an arbitrary $\hat{\mathbf{v}}(\mathbf{y})$ is chosen.

- (c) Decoding \mathbf{w} : At the third stage, the second decoder looks for the unique vector

$$\mathbf{w} \in \mathcal{C}_B(k, \hat{\mathbf{u}}(\mathbf{y}), \hat{\mathbf{v}}(\mathbf{y}))$$

such that

$$(\hat{\mathbf{u}}(\mathbf{y}), \hat{\mathbf{v}}(\mathbf{y}), \mathbf{w}, \mathbf{y}) \in \mathcal{T}_{UVWY}^{5\delta|\mathcal{X}|}. \quad (148)$$

Denote this vector by $\hat{\mathbf{w}}(\mathbf{y})$. If such a vector does not exist, or there is more than one, an arbitrary representative is chosen.

(d) At the last stage, the second reproduction vector $\hat{\mathbf{x}}_2$ is constructed, as

$$\hat{\mathbf{x}}_2 = (\hat{x}_{2,1}, \hat{x}_{2,2}, \dots, \hat{x}_{2,n}) \quad (149)$$

where

$$\hat{x}_{2,l} = f_2(\hat{w}_l(\mathbf{y}), y_l). \quad (150)$$

Analysis of the Probability of Error

The probability of error in the encoding/decoding scheme is examined next. We start by defining the error events. The error events in the encoding scheme are the following:

$$\begin{aligned} E_0 &= \{ \mathbf{x} \in (\mathcal{T}_X^\delta)^c \} \\ E_1 &= E_0^c \cap \left\{ \bigcap_{i=1}^{M_1} \{ (\mathbf{x}, \mathbf{u}_i) \notin \mathcal{T}_{XU}^{2\delta} \} \right\} \\ E_2 &= E_0^c \cap E_1^c \cap \left\{ \bigcap_{j=1}^{M_2} \{ (\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}_j) \notin \mathcal{T}_{XUV}^{3\delta} \} \right\} \\ E_3 &= \left(\bigcap_{m=0}^2 E_m^c \right) \cap \left\{ \bigcap_{k=1}^{M_3} \{ (\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}), \mathbf{w}_k) \notin \mathcal{T}_{XUVW}^{4\delta} \} \right\} \end{aligned}$$

where E_1 and E_2 correspond to the first and second stages of the first encoder, and E_3 corresponds to the second encoder. The error events of decoder 2 are

$$\begin{aligned} E_4 &= \left(\bigcap_{m=0}^3 E_m^c \right) \cap \{ (\mathbf{u}(\mathbf{x}), \mathbf{x}, \mathbf{y}) \notin \mathcal{T}_{UXY}^{3\delta} \} \\ E_5 &= \left(\bigcap_{m=0}^4 E_m^c \right) \cap \left\{ \bigcup_{\mathbf{u} \in \mathcal{A}_B(i), \mathbf{u} \neq \mathbf{u}(\mathbf{x})} \{ (\mathbf{u}, \mathbf{y}) \in \mathcal{T}_{UY}^{3\delta|\mathcal{X}|} \} \right\} \\ E_6 &= \left(\bigcap_{m=0}^5 E_m^c \right) \cap \{ (\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}), \mathbf{x}, \mathbf{y}) \notin \mathcal{T}_{UVXY}^{4\delta} \} \\ E_7 &= \left(\bigcap_{m=0}^6 E_m^c \right) \cap \left\{ \bigcup_{\mathbf{v} \in \mathcal{B}_B(j, \mathbf{u}(\mathbf{x}))} \{ (\mathbf{u}(\mathbf{x}), \mathbf{v}, \mathbf{y}) \in \mathcal{T}_{UVY}^{4\delta|\mathcal{X}|} \} \right\} \\ E_8 &= \left(\bigcap_{m=0}^7 E_m^c \right) \cap \{ (\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}), \mathbf{w}(\mathbf{x}), \mathbf{x}, \mathbf{y}) \notin \mathcal{T}_{UVWXY}^{5\delta} \} \end{aligned}$$

$$E_9 = \left(\bigcap_{m=0}^8 E_m^c \right) \cap \left\{ \bigcup_{\mathbf{w} \in \mathcal{C}_B(k, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))} \{(\mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}), \mathbf{w}, \mathbf{y}) \in \mathcal{T}_{UVWY}^{5\delta|\mathcal{X}|}\} \right\}$$

Finally, the error events for the first decoder are defined as

$$E_{10} = \left(\bigcap_{m=0}^3 E_m^c \right) \cap \{(\mathbf{u}(\mathbf{x}), \mathbf{x}, \mathbf{z}) \notin \mathcal{T}_{UXZ}^{3\delta}\}$$

$$E_{11} = \left(\bigcap_{m=0}^3 E_m^c \right) \cap E_{10}^c \cap \left\{ \bigcup_{\mathbf{u} \in \mathcal{A}_B(i)} \{(\mathbf{u}, \mathbf{z}) \in \mathcal{T}_{UZ}^{3\delta|\mathcal{X}|}\} \right\}$$

The probability of error in the encoding/decoding process, P_e , is upper bounded as

$$P_e \leq P \left\{ \bigcup_{\ell=0}^{11} E_\ell \right\} \leq \sum_{\ell=0}^{11} P(E_\ell). \quad (151)$$

Observe that if no error occurs in the encoding/decoding process (i.e., the event $\bigcap_{m=0}^{11} E_m^c$ occurs), then the following is satisfied

$$(\hat{\mathbf{W}}(\mathbf{Y}), \mathbf{X}, \mathbf{Y}) \in \mathcal{T}_{WXY}^\mu \quad (152)$$

$$(\hat{\mathbf{U}}(\mathbf{Y}), \mathbf{X}, \mathbf{Z}) \in \mathcal{T}_{UXZ}^\mu \quad (153)$$

where $\mu = 5\delta \cdot |\mathcal{U} \times \mathcal{V}|$, which means that the empirical distributions of these random vectors is close to their corresponding joint distributions. Thus, in view of (12), (13), the distortion constraints are approximately satisfied.

Thus, to prove the direct part of Theorem 1, it is enough to show that for fixed γ and sufficiently small δ , each of the terms in the sum in (151) vanishes as $n \rightarrow \infty$. The techniques for proving these limits are now classical, and follow those of [12] (see also [1] and [5]). We will focus on events E_0 to E_7 . The proof that the probabilities of the events E_8 to E_{12} vanish, follows by similar arguments.

E_0 . Clearly, $P(E_0)$ tends to zero as $n \rightarrow \infty$.

E_1 . For E_1 we have

$$\begin{aligned} P(E_1) &= P \left(\bigcap_{i=1}^{M_1} \{(\mathbf{X}, \mathbf{U}_i) \notin \mathcal{T}_{XU}^{2\delta}\} \mid \mathbf{X} \in \mathcal{T}_X^\delta \right) P(\mathbf{X} \in \mathcal{T}_X^\delta) \\ &= \sum_{\mathbf{x} \in \mathcal{T}_X^\delta} P_U \left(\bigcap_{i=1}^{M_1} \{(\mathbf{x}, \mathbf{U}_i) \notin \mathcal{T}_{XU}^{2\delta}\} \right) P_X(\mathbf{x}) \end{aligned} \quad (154)$$

By the left hand side of (131), for $\mathbf{x} \in \mathcal{T}_X^\delta$ we have

$$\begin{aligned}
& P_U \left(\bigcap_{i=1}^{M_1} \{(\mathbf{x}, \mathbf{U}_i) \notin \mathcal{T}_{XU}^{2\delta}\} \right) \\
& \leq \{1 - \exp[-nI(X; U) - n\epsilon_u]\}^{M_1} \\
& \leq \exp[-\exp(n\gamma - n\epsilon_u)]
\end{aligned} \tag{155}$$

which tends to 0 doubly exponentially fast, provided $\gamma > \epsilon_u(\delta, 2\delta)$. Substituting (155) in (154), we have

$$\lim_{n \rightarrow \infty} P(E_1) \rightarrow 0. \tag{156}$$

E_2 . Conditioned on $E_0^c \cap E_1^c$, we have $(\mathbf{x}, \mathbf{u}(\mathbf{x})) \in \mathcal{T}_{XU}^{2\delta}$, thus similarly to (154) we obtain

$$P(E_2) = \sum_{(\mathbf{x}, \mathbf{u}) \in \mathcal{T}_{XU}^{2\delta}} P_{V|U} \left(\bigcap_{j=1}^{M_2} (\mathbf{x}, \mathbf{u}, \mathbf{V}_j) \notin \mathcal{T}_{XUV}^{3\delta} \mid \mathbf{u} \right) P_{XU}(\mathbf{x}, \mathbf{u}) \tag{157}$$

By the left hand side of (132), for $(\mathbf{x}, \mathbf{u}) \in \mathcal{T}_{XU}^{2\delta}$,

$$\begin{aligned}
& P_{V|U} \left(\bigcap_{j=1}^{M_2} \{(\mathbf{x}, \mathbf{u}, \mathbf{V}_j) \notin \mathcal{T}_{XUV}^{3\delta}\} \mid \mathbf{u} \right) \\
& \leq \{1 - \exp[-nI(X; V|U) - n\eta_{v|u}]\}^{M_2} \\
& \leq \exp[-\exp(n\gamma - n\eta_{v|u})]
\end{aligned} \tag{158}$$

which vanishes provided $\gamma > \eta_{v|u}(2\delta, 3\delta)$. Hence

$$\lim_{n \rightarrow \infty} P(E_2) \rightarrow 0. \tag{159}$$

E_3 . Conditioned on $\bigcap_{m=0}^2 E_m^c$, the triplet $(\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))$ is jointly typical, and belongs to $\mathcal{T}_{XUV}^{3\delta}$. Thus the steps to show that $P(E_3)$ vanishes are similar to those leading from (157) to (159), except that the conditioning is on (U, V) , and the small parameter is 4δ instead of 3δ . Thus we have

$$\lim_{n \rightarrow \infty} P(E_3) \rightarrow 0. \tag{160}$$

provided $\gamma > \eta_{uvw}(3\delta, 4\delta)$.

E_4 . Conditioned on $\cap_{m=0}^3 E_m^c$, the pair $(\mathbf{x}, \mathbf{u}(\mathbf{x}))$ is jointly typical. Moreover, there is a Markov structure $U \circlearrowleft X \circlearrowleft Y$, where the random vector \mathbf{Y} is drawn according to $P_{Y|X}$. Consequently, by the Markov lemma [1, Lemma 14.8.1]

$$P(E_4) \leq \sum_{\mathbf{x}, \mathbf{u} \in \mathcal{T}_{XU}^{2\delta}} P_{Y|X} \left\{ (\mathbf{x}, \mathbf{u}, \mathbf{Y}) \notin \mathcal{T}_{XUY}^{3\delta} \mid \mathbf{x} \right\} P_{XU}(\mathbf{x}, \mathbf{u}) \rightarrow 0 \quad (161)$$

as $n \rightarrow 0$.

E_5 . The sequences $\{\mathbf{U}_i\}_{i=1}^{M_1}$ are drawn independent of everything. Conditioned on the intersection $\cap_{m=0}^4 E_m^c$, the vector \mathbf{y} is typical. The probability that independently drawn \mathbf{U}_i is jointly typical with \mathbf{y} is upper bounded by (see right hand side of (131))

$$P_U \left((\mathbf{U}, \mathbf{y}) \in \mathcal{T}_{UY}^{3\delta|\mathcal{X}} \right) \leq \exp[-nI(U; Y) + n\epsilon(2\delta|\mathcal{X}|, 3\delta|\mathcal{X}|)] \quad (162)$$

Since $|\mathcal{A}_B(i)| = \exp[n(I(Z; U) - 4\gamma)]$ and $I(Z; U) \leq I(Y; U)$, we have

$$\lim_{n \rightarrow \infty} P(E_5) = 0 \quad (163)$$

provided $4\gamma > \epsilon(2\delta|\mathcal{X}|, 3\delta|\mathcal{X}|)$.

E_6 . Conditioned on $\cap_{m=0}^5 E_m^c$, the triplet $(\mathbf{x}, \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}))$ is jointly typical. We have also the Markov structure $U \circlearrowleft V \circlearrowleft X \circlearrowleft Y$, thus we apply the Markov lemma, as in (161), to conclude that $P(E_6)$ vanishes as $n \rightarrow \infty$.

E_7 . The vectors $\{\mathbf{v}_j\}_{j=1}^{M_2}$ are drawn according to $P_{V|U}(\cdot | u_i(\mathbf{x}))$. Conditioned on the intersection $\cap_{m=0}^6 E_m^c$, the pair $(\mathbf{u}(\mathbf{x}), \mathbf{y})$ is typical. The probability that a vector \mathbf{V}_j drawn according to $P_{V|U}$ independently of \mathbf{y} is jointly typical with \mathbf{y} is upper bounded by (see right hand side of (132))

$$P_{V|U} \left((\mathbf{u}(\mathbf{x}), \mathbf{V}, \mathbf{y}) \in \mathcal{T}_{UVY}^{4\delta|\mathcal{X}} \right) \leq \exp[-n(I(Y; V|U) - \epsilon(3\delta|\mathcal{X}|, 4\delta|\mathcal{X}|))] \quad (164)$$

Since

$$|\mathcal{B}_B(j, \mathbf{u}(\mathbf{x}))| = \exp[n(I(Y; V|U) - 4\gamma)] \quad (165)$$

we have

$$\lim_{n \rightarrow \infty} P(E_7) = 0 \quad (166)$$

provided $4\gamma > \epsilon(3\delta|\mathcal{X}|, 4\delta|\mathcal{X}|)$.

□

6.2 Proof of Theorem 4

The direct part follows from Theorem 3, by substituting $V_{k,k} = U_k$, and choosing $V_{k,l}$, $k < l \leq K$, to be null random variables. For the converse, we will show the existence of a K -tuple of random variables (U_1, U_2, \dots, U_K) satisfying Conditions 1–3 in the definition of $\mathcal{R}_{K,i}^*(\mathbf{D})$. Let $T_i = \phi_i(X^n)$ stand for the output of the i 'th encoder, $i = 1, 2, \dots, K$. The following chains of inequalities hold

$$\begin{aligned}
 nR_1 &\geq H(T_1) \geq I(X^n; T_1 | Y^n) = H(X^n | Y^n) - H(X^n | T_1, Y^n) \\
 &= \sum_{i=1}^n \left[H(X_i | Y^n X^{i-1}) - H(X_i | T_1 Y^n X^{i-1}) \right] \\
 &= \sum_{i=1}^n \left[H(X_i | Y_i) - H(X_i | Y_i T_1 Y^{n \setminus i} X^{i-1}) \right] \\
 &= \sum_{i=1}^n I(X_i; T_1 Y^{n \setminus i} X^{i-1} | Y_i), \tag{167}
 \end{aligned}$$

$$\begin{aligned}
 n(R_2 - R_1) &\geq H(T_2 | T_1) \geq I(X^n; T_2 | T_1, Y^n) \\
 &= \sum_{i=1}^n I(X_i; T_2 | Y_i, T_1 Y^{n \setminus i} X^{i-1}), \tag{168}
 \end{aligned}$$

and, in general,

$$\begin{aligned}
 n(R_k - R_{k-1}) &\geq H(T_k | T_{k-1}, \dots, T_1) \geq I(X^n; T_k | T_{k-1}, \dots, T_1) \\
 &= \sum_{i=1}^n I(X_i; T_k | Y_i, T_{k-1} \dots T_1 Y^{n \setminus i} X^{i-1}). \tag{169}
 \end{aligned}$$

Define the random variables

$$U_{1,i} \triangleq T_1 Y^{n \setminus i} X^{i-1} \tag{170}$$

$$U_{k,i} \triangleq T_k U_{k-1,i}, \quad k = 2, 3, \dots, K. \tag{171}$$

With these definitions, the proof proceeds along the lines leading from (101) to (111). The details, being straightforward, are omitted. \square

References

- [1] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.

- [2] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Academic Press, London, 1981.
- [3] W. H. R. Equitz, *Successive refinement of information*, Ph.D. dissertation, Stanford University, June 1989.
- [4] W. H. R. Equitz and T. M. Cover, "Successive refinement of information," *IEEE Trans. Inform. Theory*, vol. 37, pp. 269-274, Mar. 1991.
- [5] C. Heegard and T. Berger, "Rate distortion when side information may be absent," *IEEE Trans. Inform. Theory*, vol. IT-31, pp. 727-734, Nov. 1985.
- [6] A. Kaspi, "Rate-distortion function when side-information may be present at the decoder," *IEEE Trans. Inform. Theory*, vol. 40, pp. 2031-2034, Nov. 1994.
- [7] V. N. Koshelev, "Hierarchical coding of discrete sources," *Probl. Peredachi Inform.*, vol. 16, no. 3, pp. 31-49, 1980. English translation: vol. 16, pp. 186-203, 1980.
- [8] V. N. Koshelev, "On the divisibility of discrete sources with an additive single-letter distortion measure," *Probl. Peredachi Inform.*, vol. 30, no. 1, pp. 31-50, 1994. English translation: vol. 30, no. 1, pp. 27-43, 1994.
- [9] B. Rimoldi, "Successive refinement of information: Characterization of achievable rates," *IEEE Trans. Inform. Theory*, vol 40, pp. 253-259, Jan. 1994.
- [10] S. Shamai and S. Verdú, "Capacity of channels with uncoded side information," *Europ. Trans. Telecommun.*, vol. 6, no. 5, pp. 587-600, Sept.-Oct. 1995.
- [11] S. Shamai, S. Verdú, and R. Zamir, "Systematic lossy source/channel coding," *IEEE Trans. Inform. Theory*, vol. 44, no. 2, pp. 564-579, March 1998.
- [12] A. D. Wyner and J. Ziv, "The rate-distortion function for source coding with side information at the decoder," *IEEE Trans. Inform. Theory*, vol. IT-22, pp. 1-10, Jan. 1976.