

# Distributed Information-Theoretic Biclustering

Georg Pichler\*, Pablo Piantanida†, and Gerald Matz\*

\*Institute of Telecommunications, Vienna University of Technology, Vienna, Austria

†Laboratoire des Signaux et Systèmes (L2S), CentraleSupélec-CNRS-Université Paris-Sud, Gif-sur-Yvette, France  
Email: {georg.pichler,gerald.matz}@nt.tuwien.ac.at; pablo.piantanida@centralesupelec.fr

**Abstract**—This paper investigates the problem of distributed biclustering of memoryless sources and extends previous work [1] to the general case with more than two sources. Given a set of distributed stationary memoryless sources, the encoders’ goal is to find rate-limited representations of these sources such that the mutual information between two selected subsets of descriptions (each of them generated by distinct encoders) is maximized. This formulation is fundamentally different from conventional distributed source coding problems since here redundancy among descriptions should actually be maximally preserved. We derive non-trivial outer and inner bounds to the achievable region for this problem and further connect them to the CEO problem under logarithmic loss distortion. Since information-theoretic biclustering is closely related to distributed hypothesis testing against independence, our results are also expected to apply to that problem.

## I. INTRODUCTION

Clustering unstructured data is desirable to understand the nature of unstructured data with the goal of making information more accessible to humans. Clustering algorithms have been very successfully used in wide range of applications in areas like life sciences (gene expression analysis [2], segmentation of PET images [3]) and marketing research [4].

In this paper, we study an information-theoretic formulation of the biclustering problem that extends our recent work [1] to the case of multiple (usually dependent) sources. Each source is assumed to be stationary and memoryless and is observed by a distinct encoder. The encoders aim at extracting lossy (rate-limited) descriptions of the associated source such that these descriptions are maximally informative about each other. In this setting, we investigate the optimal tradeoff between *relevance*, the normalized multi-letter mutual information between the descriptions, and *complexity*, the encoding rates. We refer to this scenario as the distributed multi-terminal *information-theoretic biclustering* (ITB) problem and derive an inner bound to the sets of achievable rates of relevance and complexity.

The problem at hand is fundamentally different from conventional distributed source coding problems, which aim at discarding redundant information while guaranteeing correct decoding. In contrast to the ITB problem, where redundant information should actually be maximally preserved. Furthermore, the ITB problem is closely related to testing against independence with multi-terminal data compression [5], for which a general solution remains elusive [6].

While the outer bound follows directly from the case of two sources, the straightforward generalization of the achievability result [1, Theorem 6] yields a very poor inner bound. This bound can be improved by correctly accounting for the possibility of encoders that act as helpers, which can be achieved using a binning strategy, that is, a variant of Berger-Tung coding. Nevertheless, the inner bound cannot be tight, as the well-known Körner-Marton problem [7] becomes a special case of the considered problem when considering more than two sources. Thus, in general Berger-Tung coding is suboptimal.

We further investigate the CEO problem under a mutual information constraint. This is another special case of ITB, and we show that it becomes equivalent to classical multiterminal lossy source coding under logarithmic loss distortion. By leveraging this equivalence, we obtain tight bounds for this special case using methods from [8].

## Notation

We denote random quantities and their realizations by capital, sans-serif and lowercase letters, respectively. Furthermore, vectors are indicated by bold-face type and sets by calligraphic type. We use subscript to denote slices of vectors, i.e.,  $\mathbf{x}_{\mathcal{A}} \triangleq (\mathbf{x}_i)_{i \in \mathcal{A}}$  as well as the usual notation  $\mathbf{x}_i^j \triangleq \mathbf{x}_{\{i, \dots, j\}}$ ,  $\mathbf{x}^j \triangleq \mathbf{x}_j^j$ . If a vector is already carrying a subscript, it will be separated by a comma, e.g.,  $\mathbf{x}_{3,1}^5 = (\mathbf{x}_3)_1^5 = (\mathbf{x}_3)^5$ . Random variables are assumed to be supported on finite sets. We use the same letter for the random variable and for its support set, e.g.,  $\mathcal{Y}$  takes values in  $\mathcal{Y}$  and  $X_3$  takes values in  $\mathcal{X}_3$ . For a set  $\mathcal{X}$ , let  $\mathcal{X}^n$  denote the  $n$ -th Cartesian power of  $\mathcal{X}$ . Given a random variable  $X$ , we write  $p_X \in \mathcal{P}(\mathcal{X})$  for its probability mass function (pmf), where  $\mathcal{P}(\mathcal{X})$  is the set of all pmfs on  $\mathcal{X}$ . We write  $X \sim p$  to indicate that  $X$  is distributed according to  $p \in \mathcal{P}(\mathcal{X})$ . We use the notation of [9, Chapter 2] for information-theoretic quantities, however, all logarithms in this paper are to base  $e$  and information theoretic quantities are measured in nats. The notation  $h_0(p) \triangleq -p \log p - (1-p) \log(1-p)$  is used for the binary entropy function and  $X \rightarrow Y \leftarrow Z$  indicates that  $X$ ,  $Y$ , and  $Z$  form a Markov chain in this order. For convenience, we define  $\mathcal{K} \triangleq \{1, 2, \dots, K\}$  and we use the symbol  $x_{2\kappa, 2\kappa}$  for a tuple, indexed by disjoint unordered pairs  $\{\mathcal{A}, \mathcal{B}\}$  of nonempty sets  $\mathcal{A}, \mathcal{B} \subset \mathcal{K}$  and write  $x_{\mathcal{A}, \mathcal{B}}$  for its components. For a random variable  $X$ , we denote the set of (strongly)  $\delta$ -typical  $n$ -sequences [10, Section 2.4] as  $\mathcal{T}_{[X]^\delta}^n$ . Regarding the method of types, which is used heavily in Section IV, the reader is referred to [11].

## II. PROBLEM DEFINITION AND MAIN RESULTS

### A. Problem Statement

Let  $X_{\mathcal{K}}$  be  $K$  random variables. The random vectors  $\mathbf{X}_{\mathcal{K}}$  consist of  $n$  i.i.d. copies of  $X_{\mathcal{K}}$ . For  $n \in \mathbb{N}$  and  $R_{\mathcal{K}} \in \mathbb{R}^K$ , an  $(n, R_{\mathcal{K}})$  code  $f_{\mathcal{K}}$  consists of  $K$  functions  $f_k: \mathcal{X}_k^n \rightarrow \mathcal{M}_k$ , where  $\mathcal{M}_k$  is an arbitrary finite set with  $\log |\mathcal{M}_k| \leq nR_k$  for each  $k \in \mathcal{K}$ .

**Definition 1** (Relevance). *Consider an  $(n, R_{\mathcal{K}})$  code  $f_{\mathcal{K}}$  and let  $W_k \triangleq f_k(\mathbf{X}_k)$ . For any pair of disjoint nonempty sets  $\mathcal{A}, \mathcal{B} \subset \mathcal{K}$ , we define the co-information of  $f_{\mathcal{A}}$  and  $f_{\mathcal{B}}$  as*

$$\Theta(f_{\mathcal{A}}; f_{\mathcal{B}}) \triangleq \frac{1}{n} \mathbb{I}(W_{\mathcal{A}}; W_{\mathcal{B}}).$$

**Definition 2** (Achievability and relevance-rates region). A point  $(\mu_{2\mathcal{K}, 2\mathcal{K}}, R_{\mathcal{K}})$  is achievable if there exists an  $(n, R_{\mathcal{K}})$  code  $f_{\mathcal{K}}$  for some  $n \in \mathbb{N}$  such that

$$\Theta(f_{\mathcal{A}}; f_{\mathcal{B}}) \geq \mu_{\mathcal{A}, \mathcal{B}}, \quad \forall \mathcal{A}, \mathcal{B} \subset \mathcal{K} \text{ disjoint, nonempty.}$$

The achievable relevance-rates region  $\overline{\mathcal{R}}$  is the closure of the set  $\mathcal{R}$  of achievable points.

### B. Main Results

**Theorem 3** (Outer bound). We have  $\mathcal{R} \subseteq \mathcal{R}_o$ , where the outer bound  $\mathcal{R}_o$  is the set of all points  $(\mu_{2\mathcal{K}, 2\mathcal{K}}, R_{\mathcal{K}})$  such that there exist random variables  $U_{\mathcal{K}}$  with  $U_{\mathcal{A}} \ominus X_{\mathcal{A}} \ominus X_{\mathcal{K}}$  for any  $\mathcal{A} \subseteq \mathcal{K}$  and

$$\sum_{k \in \mathcal{A}} R_k \geq I(U_{\mathcal{A}}; X_{\mathcal{K}} | U_{\mathcal{C}}) \quad \forall \mathcal{C} \subseteq \mathcal{K},$$

$$\mu_{\mathcal{A}, \mathcal{B}} \leq I(U_{\mathcal{A}}; X_{\mathcal{A}}) + I(U_{\mathcal{B}}; X_{\mathcal{B}}) - I(U_{\mathcal{A}} U_{\mathcal{B}}; X_{\mathcal{A}} X_{\mathcal{B}}),$$

for all  $k \in \mathcal{K}$  and all disjoint nonempty sets  $\mathcal{A}, \mathcal{B} \subset \mathcal{K}$ .

*Proof:* The proof of this theorem follows from standard information-theoretic arguments using the auxiliary random variable identification  $U_{k,i} \triangleq (f(\mathbf{X}_k), \mathbf{X}_{k,1}^{i-1})$ . ■

**Definition 4.** A point  $(\mu_{2\mathcal{K}, 2\mathcal{K}}, R_{\mathcal{K}})$  is in the region  $\mathcal{R}_i$  if there exist random variables  $U_{\mathcal{K}}$  satisfying  $U_k \ominus X_k \ominus (X_{\mathcal{K} \setminus k}, U_{\mathcal{K} \setminus k})$  for all  $k \in \mathcal{K}$ , and if for every pair of disjoint nonempty index sets  $\mathcal{A}, \mathcal{B} \subset \mathcal{K}$  there exist subsets  $\mathcal{A}_b \subseteq \mathcal{A}_a \subseteq \mathcal{A}$  and  $\mathcal{B}_b \subseteq \mathcal{B}_a \subseteq \mathcal{B}$  such that

$$\sum_{k \in \mathcal{A}'} R_k \geq I(X_{\mathcal{A}'}; U_{\mathcal{A}'} | U_{\mathcal{A}_a \setminus \mathcal{A}'}), \quad (1)$$

$$\sum_{k \in \mathcal{B}'} R_k \geq I(X_{\mathcal{B}'}; U_{\mathcal{B}'} | U_{\mathcal{B}_a \setminus \mathcal{B}'}), \quad (2)$$

$$\mu_{\mathcal{A}, \mathcal{B}} \leq I(U_{\mathcal{A}_b}; U_{\mathcal{B}_b}), \quad (3)$$

for all subsets  $\mathcal{A}' \subseteq \mathcal{A}_a$  with  $\mathcal{A}' \cap \mathcal{A}_b \neq \emptyset$  and  $\mathcal{B}' \subseteq \mathcal{B}_a$  with  $\mathcal{B}' \cap \mathcal{B}_b \neq \emptyset$ .

We use typicality coding and binning to show that  $\mathcal{R}_i$  is indeed achievable. The conditions (1) and (2) ensure that  $U_{\mathcal{A}_b}$  and  $U_{\mathcal{B}_b}$  can be correctly decoded from the output of the encoders  $\mathcal{A}_a$  and  $\mathcal{B}_a$ , respectively. By (3),  $U_{\mathcal{A}_b}$  and  $U_{\mathcal{B}_b}$  are enough to ensure that  $\mu_{\mathcal{A}, \mathcal{B}}$  is achievable. Intuitively, one can say that the encoders  $\mathcal{A}_a \setminus \mathcal{A}_b$  and  $\mathcal{B}_a \setminus \mathcal{B}_b$  act as helpers for decoding  $U_{\mathcal{A}_b}$  and  $U_{\mathcal{B}_b}$ , respectively. The special case  $\mathcal{A}_b = \mathcal{A}$ ,  $\mathcal{B}_b = \mathcal{B}$  for every disjoint pair of nonempty sets  $\mathcal{A}, \mathcal{B} \subset \mathcal{K}$  corresponds to no binning at all, as (1) and (2) then imply  $R_k \geq I(X_k; U_k)$  for all  $k \in \mathcal{K}$ .

**Theorem 5** (Inner bound).  $\mathcal{R}_i \subseteq \overline{\mathcal{R}}$ .

*Proof:* See Section IV. ■

This achievability result cannot be shown by merely applying Berger-Tung coding. In contrast to averaged per-letter distortion, joint typicality alone is not sufficient to ensure that relevance is high enough. Therefore, the proof of this achievable region utilizes more sophisticated tools developed for analyzing hypothesis testing problems [5].

*Remark.* Note that the inner bound in Theorem 5 cannot be tight in general as Definition 2 contains the Körner-Marton problem [7] as a special case. Choose for  $K = 3$ ,  $X_1 \sim \mathcal{B}(\frac{1}{2})$  and  $X_3 \sim \mathcal{B}(p)$  with  $p \in (0, 1)$  and  $p \neq \frac{1}{2}$ . Then define  $X_2 \triangleq X_1 \oplus X_3$ . The point  $(\mu_{2\mathcal{A}, 2\mathcal{B}}, R_{\mathcal{K}})$  where  $R_3 = \log(2)$ ,  $R_1 = R_2 = H(X_3) = h_0(p)$ , and  $\mu_{\mathcal{A}, \mathcal{B}} = 0$  except

for  $\mu_{\{1,2\}, \{3\}} = H(X_3) = h_0(p)$  is achievable [7, Theorem 1]. However, the quantize-and-bin scheme cannot achieve this point [7, Proposition 1].

$\mathcal{R}_i$  is not convex in general, thus Theorem 5 can be strengthened to  $\text{conv}(\mathcal{R}_i) \subseteq \overline{\mathcal{R}}$ . However, characterizing  $\text{conv}(\mathcal{R}_i)$  using a time-sharing random variable is tedious, due to the choice of index sets  $\mathcal{A}_a, \mathcal{A}_b, \mathcal{B}_a$ , and  $\mathcal{B}_b$  in Definition 4.

The following result establishes cardinality bounds that render  $\mathcal{R}_i$  computable.

**Proposition 6.** The region  $\mathcal{R}_i$  remains the same if the cardinality bound  $|\mathcal{U}_k| \leq |\mathcal{X}_k| + 4^K$  is imposed for every  $k \in \mathcal{K}$ .

*Proof:* The proof is a straightforward application of the support lemma [10, Appendix C]. ■

### III. A SPECIAL CASE: THE CEO PROBLEM

In this section we will analyze a special case of the biclustering problem, a variant of the CEO problem [12], where the quality is measured in terms of mutual information (MI) instead of a distortion criterion. Consider random variables  $(X_{\mathcal{J}}, Y_{\mathcal{L}})$  with  $\mathcal{J} = \{1, 2, \dots, J\}$  and  $\mathcal{L} = \{1, 2, \dots, L\}$ . Slightly abusing notation, we write  $x_{2\mathcal{J}, 2\mathcal{L}}$  for a tuple indexed by a pair of nonempty sets  $(\mathcal{A}, \mathcal{B})$  where  $\mathcal{A} \subseteq \mathcal{J}$  and  $\mathcal{B} \subseteq \mathcal{L}$ .

**Definition 7.** A point  $(\mu_{2\mathcal{J}, 2\mathcal{L}}, R_{\mathcal{J}})$  is MI-achievable, if for some  $n \in \mathbb{N}$  there exists an  $(n, R_{\mathcal{J}})$  code  $f_{\mathcal{J}}$  for  $X_{\mathcal{J}}$  with  $U_j \triangleq f_j(\mathbf{X}_j)$  such that for all  $\mathcal{A} \subseteq \mathcal{J}$  and  $\mathcal{B} \subseteq \mathcal{L}$

$$\frac{1}{n} I(U_{\mathcal{A}}; \mathbf{Y}_{\mathcal{B}}) \geq \mu_{\mathcal{A}, \mathcal{B}}.$$

Let  $\mathcal{R}_{\text{MI}}$  denote the set of all MI-achievable points.

We can obtain  $\mathcal{R}_{\text{MI}}$  as a special case from  $\mathcal{R}$ . Consider  $X_{\mathcal{K}} = (X_{\mathcal{J}}, Y_{\mathcal{L}})$  with  $J + L = K$ . We want to find achievable points  $(\mu_{2\mathcal{K}, 2\mathcal{K}}, R_{\mathcal{K}}) \in \mathcal{R}$  that correspond to the CEO problem. Let  $Y_{\mathcal{L}}$  be transmitted without compression, i.e.,  $R_{J+l} = \log|\mathcal{Y}_l|$  for  $l \in \mathcal{L}$ . Furthermore, we are only interested in the information the encodings of  $\mathbf{X}_{\mathcal{J}}$  provide about  $\mathbf{Y}_{\mathcal{L}}$ , i.e.,  $\mu_{\mathcal{A}, \mathcal{B}} \neq 0$  only if  $\mathcal{A} \subseteq \mathcal{J}$  and  $\mathcal{B} \subseteq \mathcal{L}$ .

We want to show that the above CEO problem is equivalent to the logarithmic loss (log-loss) distortion approach from [8]. For  $\mathcal{A} \subseteq \mathcal{J}$  and  $\mathcal{B} \subseteq \mathcal{L}$ , we define a decoding function  $g_{\mathcal{A}, \mathcal{B}}: \mathcal{M}_{\mathcal{A}} \rightarrow \mathcal{P}(\mathcal{Y}_{\mathcal{B}}^n)$ , which produces a probabilistic estimate of  $\mathbf{Y}_{\mathcal{B}}$  given the output of the encoders  $\mathcal{A}$ . We use log-loss fidelity  $\zeta_{\mathcal{B}}, \mathcal{B} \subseteq \mathcal{L}$ , defined as the negative log-loss distortion  $d_{\mathcal{B}}(\mathbf{y}_{\mathcal{B}}, p) \triangleq -\frac{1}{n} \log p(\mathbf{y}_{\mathcal{B}})$  on  $\mathcal{Y}_{\mathcal{B}}^n$  augmented by entropy, i.e.,

$$\zeta_{\mathcal{B}}: \mathcal{Y}_{\mathcal{B}}^n \times \mathcal{P}(\mathcal{Y}_{\mathcal{B}}^n) \rightarrow \mathbb{R},$$

$$(\mathbf{y}_{\mathcal{B}}, p) \mapsto H(\mathbf{Y}_{\mathcal{B}}) - d_{\mathcal{B}}(\mathbf{y}_{\mathcal{B}}, p).$$

**Definition 8** (Rate-distortion region for log-loss distortion). A point  $(\mu_{2\mathcal{J}, 2\mathcal{L}}, R_{\mathcal{J}})$  is log-loss achievable if for some  $n \in \mathbb{N}$  there exists an  $(n, R_{\mathcal{J}})$ -code  $f_{\mathcal{J}}$  yielding  $W_j \triangleq f_j(\mathbf{X}_j)$  and decoding functions  $g_{\mathcal{A}, \mathcal{B}}: \mathcal{M}_{\mathcal{A}} \rightarrow \mathcal{P}(\mathcal{Y}_{\mathcal{B}}^n)$  for all nonempty  $\mathcal{A} \subseteq \mathcal{J}$  and  $\mathcal{B} \subseteq \mathcal{L}$  such that

$$\mathbb{E}[\zeta_{\mathcal{B}}(\mathbf{Y}_{\mathcal{B}}, g_{\mathcal{A}, \mathcal{B}}(W_{\mathcal{A}}))] \geq \mu_{\mathcal{A}, \mathcal{B}}.$$

Let  $\mathcal{R}_{\text{LL}}$  be the set of all log-loss achievable points.

To show the equivalence with the ITB problem, we state an auxiliary lemma similar to [8, Lemma 1].

**Lemma 9.** For any decoding function  $g_{\mathcal{A},\mathcal{B}}$  and code  $f_{\mathcal{J}}$ ,

$$\mathbb{E}[\zeta_{\mathcal{B}}(\mathbf{Y}_{\mathcal{B}}, g_{\mathcal{A},\mathcal{B}}(\mathbf{U}_{\mathcal{A}}))] \leq \frac{1}{n} \mathbb{I}(\mathbf{Y}_{\mathcal{B}}; \mathbf{U}_{\mathcal{A}}),$$

with equality if and only if  $g_{\mathcal{A},\mathcal{B}}(\mathbf{u}_{\mathcal{A}}) = \text{p}_{\mathbf{Y}_{\mathcal{B}}|\mathbf{U}_{\mathcal{A}}}(\cdot | \mathbf{u}_{\mathcal{A}})$ .

*Proof:* The proof follows with suitable modifications directly from the proof of [8, Lemma 1]. ■

As a corollary of Lemma 9, we have equality between the two achievable regions.

**Corollary 10** (Equivalence of MI and log-loss).  $\mathcal{R}_{LL} = \mathcal{R}_{MI}$ .

*Proof:* Immediate from Lemma 9. ■

Corollary 10 implies that the results in [8] directly apply to the CEO problem with mutual information constraint. For instance, the CEO problem with  $J$  encoders under logarithmic loss distortion [8, Appendix B] can be obtained in Definition 8 by setting  $L = 1$  (i.e., the CEO is only interested in  $\mathbf{Y}_1$ ) and  $\mu_{\mathcal{A},1} = 0$ , whenever  $\mathcal{A} \neq \mathcal{J}$  (i.e., the CEO listens to all her agents). Then the region  $\mathcal{R}_{MI} = \mathcal{R}_{LL}$  collapses to the set of all  $(\mu, R_{\mathcal{J}})$  such that there exists an  $(n, R_{\mathcal{J}})$ -code  $f_{\mathcal{J}}$  with  $\mathbb{I}(\mathbf{U}_{\mathcal{J}}; \mathbf{Y}_1) \geq n\mu$ , where  $\mathbf{U}_j \triangleq f_j(\mathbf{X}_j)$ . The resulting achievable region  $\mathcal{R}_{MI}^i \subseteq \mathcal{R}_{MI}$ , that can be obtained from  $\mathcal{R}_i$  is shown to be tight by [8, Lemma 5] if  $\mathcal{X}_{\mathcal{J}}$  are mutually independent given  $\mathbf{Y}_1$ . A similar argument can be made for the multiterminal source coding problem under logarithmic loss, introduced in [8, Section II], where the inner bound, obtained from  $\mathcal{R}_i$ , is also tight.

#### IV. PROOF OF THEOREM 5

This proof extends ideas in [5, Section VI] to more than two sources and incorporates a binning strategy.

**Lemma 11** (Existence of a code). Let  $\varepsilon > 0$ ,  $\mathbf{U}_k \text{---} \mathbf{X}_k$ ,  $\text{---} (\mathbf{X}_{\mathcal{K} \setminus k}, \mathbf{U}_{\mathcal{K} \setminus k})$  for all  $k \in \mathcal{K}$ , and  $R_{\mathcal{K}} \in \mathbb{R}_+^{\mathcal{K}}$ . Then, for some suitably small  $\delta > 0$  and suitably large  $n \in \mathbb{N}$ , we can obtain an  $(n, R_{\mathcal{K}} + \varepsilon)$  code  $f_{\mathcal{K}}$  and decoding functions  $g_{\mathcal{A}_a, \mathcal{A}_b} : \mathcal{M}_{\mathcal{A}_a} \rightarrow \mathcal{U}_{\mathcal{A}_b}^n$  for each pair of nonempty index sets  $\mathcal{A}_b \subseteq \mathcal{A}_a \subseteq \mathcal{K}$ , such that the following two properties hold:

- 1) Let  $\mathcal{A}_a, \mathcal{A}_b, \mathcal{B}_a, \mathcal{B}_b \subseteq \mathcal{K}$  be arbitrary nonempty subsets of indices with  $\mathcal{A}_b \subseteq \mathcal{A}_a$ ,  $\mathcal{B}_b \subseteq \mathcal{B}_a$ , and  $\mathcal{A}_a \cap \mathcal{B}_a = \emptyset$ . If (1) and (2) hold, then, using  $\mathbf{W}_k \triangleq f_k(\mathbf{X}_k)$  and the abbreviations  $\hat{\mathbf{U}}_{\mathcal{A}_b} \triangleq g_{\mathcal{A}_a, \mathcal{A}_b}(\mathbf{W}_{\mathcal{A}_a})$  and  $\hat{\mathbf{U}}_{\mathcal{B}_b} \triangleq g_{\mathcal{B}_a, \mathcal{B}_b}(\mathbf{W}_{\mathcal{B}_a})$ ,

$$\mathbb{P}\left\{(\hat{\mathbf{U}}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \hat{\mathbf{U}}_{\mathcal{B}_b}) \notin \mathcal{T}_{[\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b}] \delta}^n\right\} \leq \varepsilon. \quad (4)$$

- 2) For any pair  $\mathcal{A}_a, \mathcal{B}_a \subseteq \mathcal{K}$  of disjoint nonempty index sets,

$$\left| [g_{\mathcal{A}_a, \mathcal{A}_b}(\mathcal{M}_{\mathcal{A}_a}) \times g_{\mathcal{B}_a, \mathcal{B}_b}(\mathcal{M}_{\mathcal{B}_a})] \cap \mathcal{T}_{[\mathbf{U}_{\mathcal{A}_b}, \mathbf{U}_{\mathcal{B}_b}] \delta}^n \right| \leq \exp(n(\mathbb{I}(\mathbf{U}_{\mathcal{A}_b}, \mathbf{U}_{\mathcal{B}_b}); \mathbf{X}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{B}_b}) + \varepsilon) \quad (5)$$

for any  $\mathcal{A}_b \subseteq \mathcal{A}_a$  and  $\mathcal{B}_b \subseteq \mathcal{B}_a$ .

*Proof:* See Appendix A. ■

We will further need the following set of random variables.

**Definition 12.** For random variables  $(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2)$  and  $\delta \geq 0$ , define the set of random variables  $\mathcal{L}_{\delta}(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2)$ , containing  $(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2)$  as

$$\mathcal{L}_{\delta}(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2) \triangleq \left\{ (\tilde{\mathbf{U}}_1, \tilde{\mathbf{X}}_1, \tilde{\mathbf{X}}_2, \tilde{\mathbf{U}}_2) : \forall (u_1, x_1, x_2, u_2) : \left| \text{p}_{\tilde{\mathbf{U}}_1, \tilde{\mathbf{X}}_1}(u_1, x_1) - \text{p}_{\mathbf{U}_1, \mathbf{X}_1}(u_1, x_1) \right| \leq \delta \right\}$$

$$\left| \text{p}_{\tilde{\mathbf{U}}_2, \tilde{\mathbf{X}}_2}(u_2, x_2) - \text{p}_{\mathbf{U}_2, \mathbf{X}_2}(u_2, x_2) \right| \leq \delta \text{p}_{\mathbf{U}_2, \mathbf{X}_2}(u_2, x_2), \left| \text{p}_{\tilde{\mathbf{U}}_1, \tilde{\mathbf{U}}_2}(u_1, u_2) - \text{p}_{\mathbf{U}_1, \mathbf{U}_2}(u_1, u_2) \right| \leq \delta \text{p}_{\mathbf{U}_1, \mathbf{U}_2}(u_1, u_2) \Big\}.$$

Note that  $\mathcal{L}_{\delta}(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2) \subseteq \mathcal{L}_{\delta'}(\mathbf{U}_1, \mathbf{X}_1, \mathbf{X}_2, \mathbf{U}_2)$  for  $\delta \leq \delta'$ . Let  $(\mu_{2\mathcal{K}}, \mu_{2\mathcal{K}}, R_{\mathcal{K}}) \in \mathcal{R}_i$  and choose  $\mathbf{U}_{\mathcal{K}}$  as given in Definition 4. Fix  $\varepsilon > 0$  and apply Lemma 11 to obtain encoding functions  $f_{\mathcal{K}}$  and decoding functions  $g_{\mathcal{A}_a, \mathcal{A}_b}$  for each pair of nonempty index sets  $\mathcal{A}_b \subseteq \mathcal{A}_a \subseteq \mathcal{K}$ . For an arbitrary pair of disjoint nonempty index sets  $\mathcal{A}, \mathcal{B} \subseteq \mathcal{K}$ , find the subsets  $\mathcal{A}_b \subseteq \mathcal{A}_a \subseteq \mathcal{A}$  and  $\mathcal{B}_b \subseteq \mathcal{B}_a \subseteq \mathcal{B}$ , such that (1) to (3) hold. If either  $\mathcal{A}_b = \emptyset$  or  $\mathcal{B}_b = \emptyset$ , then  $\mu_{\mathcal{A}, \mathcal{B}} \leq 0$ , which is achieved by any code. We thus assume  $\mathcal{A}_b, \mathcal{B}_b \neq \emptyset$ . Define the functions  $h_1 \triangleq g_{\mathcal{A}_a, \mathcal{A}_b} \circ f_{\mathcal{A}_a}$  and  $h_2 \triangleq g_{\mathcal{B}_a, \mathcal{B}_b} \circ f_{\mathcal{B}_a}$ . We want to analyze  $\Theta(f_{\mathcal{A}}; f_{\mathcal{B}})$ . To this end, define  $\mathcal{D}_1 \triangleq h_1(\mathcal{X}_{\mathcal{A}_a}^n)$ . Naturally we can partition  $\mathcal{X}_{\mathcal{A}_a}^n$  as  $\mathcal{X}_{\mathcal{A}_a}^n = \bigcup_{\mathbf{u}_{\mathcal{A}_b} \in \mathcal{D}_1} h_1^{-1}(\mathbf{u}_{\mathcal{A}_b})$ . We may assume without loss of generality that  $h_1^{-1}(\mathbf{u}_{\mathcal{A}_b}) \subseteq \mathcal{T}_{[\mathbf{X}_{\mathcal{A}_a}, \mathbf{U}_{\mathcal{A}_b}] \delta}^n(\mathbf{u}_{\mathcal{A}_b})$  whenever  $\mathbf{u}_{\mathcal{A}_b} \in \mathcal{D}_1 \cap \mathcal{T}_{[\mathbf{U}_{\mathcal{A}_b}] \delta}^n$  as this does not interfere with the properties of the code. Defining  $\mathcal{D}_2$  in the same manner, we set  $\mathcal{F} \triangleq (\mathcal{D}_1 \times \mathcal{D}_2) \cap \mathcal{T}_{[\mathbf{U}_{\mathcal{A}_b}, \mathbf{U}_{\mathcal{B}_b}] \delta}^n$ . Using  $\hat{\mathbf{U}}_1 \triangleq h_1(\mathbf{X}_{\mathcal{A}_a})$  and  $\hat{\mathbf{U}}_2 \triangleq h_2(\mathbf{X}_{\mathcal{B}_a})$  we have

$$\begin{aligned} n \cdot \Theta(f_{\mathcal{A}}(\mathbf{X}_{\mathcal{A}}); f_{\mathcal{B}}(\mathbf{X}_{\mathcal{B}})) & \stackrel{(a)}{\geq} n \cdot \Theta(f_{\mathcal{A}_a}(\mathbf{X}_{\mathcal{A}_a}); f_{\mathcal{B}_a}(\mathbf{X}_{\mathcal{B}_a})) \stackrel{(b)}{\geq} n \cdot \Theta(\hat{\mathbf{U}}_1; \hat{\mathbf{U}}_2) \\ & \stackrel{(c)}{\geq} \mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}\} \log \frac{\mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}\}}{\mathbb{P}\{(\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2) \in \mathcal{F}\}} \\ & \quad + \mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}^c\} \log \frac{\mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}^c\}}{\mathbb{P}\{(\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2) \in \mathcal{F}^c\}} \\ & \geq -\text{h}_0(\mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}\}) \\ & \quad - \mathbb{P}\{(\hat{\mathbf{U}}_1, \hat{\mathbf{U}}_2) \in \mathcal{F}\} \log \mathbb{P}\{(\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2) \in \mathcal{F}\} \\ & \stackrel{(d)}{\geq} -\log(2) - (1 - \varepsilon) \log \mathbb{P}\{(\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2) \in \mathcal{F}\}, \quad (6) \end{aligned}$$

where (a) and (b) follow from the data processing inequality [9, Theorem 2.8.1], (c) is a consequence of the log-sum inequality [9, Theorem 2.7.1], and (d) follows from (4). We further used  $\bar{\mathbf{U}}_1 \triangleq h_1(\bar{\mathbf{X}}_{\mathcal{A}_a})$ ,  $\bar{\mathbf{U}}_2 \triangleq h_2(\bar{\mathbf{X}}_{\mathcal{B}_a})$ , where  $(\bar{\mathbf{X}}_{\mathcal{A}_a}, \bar{\mathbf{X}}_{\mathcal{B}_a})$  are i.i.d. copies of  $(\bar{\mathbf{X}}_{\mathcal{A}_a}, \bar{\mathbf{X}}_{\mathcal{B}_a}) \sim \text{p}_{\mathbf{X}_{\mathcal{A}_a}} \text{p}_{\mathbf{X}_{\mathcal{B}_a}}$ .

For each  $\mathbf{u}_{\mathcal{A}_b} \in \mathcal{D}_1$  and  $\mathbf{u}_{\mathcal{B}_b} \in \mathcal{D}_2$ , define:

$$\mathcal{S}(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \triangleq \{\mathbf{u}_{\mathcal{A}_b}\} \times h_1^{-1}(\mathbf{u}_{\mathcal{A}_b}) \times h_2^{-1}(\mathbf{u}_{\mathcal{B}_b}) \times \{\mathbf{u}_{\mathcal{B}_b}\}$$

and

$$\mathfrak{S} \triangleq \bigcup_{(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \in \mathcal{F}} \mathcal{S}(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}).$$

Let  $(\hat{\mathbf{U}}_{\mathcal{A}_b}, \hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a}, \hat{\mathbf{U}}_{\mathcal{B}_b})$  be the type variables corresponding to some fixed  $(\hat{\mathbf{u}}_{\mathcal{A}_b}, \hat{\mathbf{x}}_{\mathcal{A}_a}, \hat{\mathbf{x}}_{\mathcal{B}_a}, \hat{\mathbf{u}}_{\mathcal{B}_b}) \in \mathfrak{S}$ . From [11, Lemma 2.6], we know

$$\mathbb{P}\{\bar{\mathbf{X}}_{\mathcal{A}_a} = \hat{\mathbf{x}}_{\mathcal{A}_a}, \bar{\mathbf{X}}_{\mathcal{B}_a} = \hat{\mathbf{x}}_{\mathcal{B}_a}\} = \exp\left[-n\left(\mathbb{H}(\hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a}) + \text{D}_{\text{KL}}(\hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a} \parallel \bar{\mathbf{X}}_{\mathcal{A}_a}, \bar{\mathbf{X}}_{\mathcal{B}_a})\right)\right]. \quad (7)$$

Letting  $K(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b})$  denote the number of elements in  $\mathcal{S}(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b})$  with type  $(\hat{\mathbf{U}}_{\mathcal{A}_b}, \hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a}, \hat{\mathbf{U}}_{\mathcal{B}_b})$ ,

$$K(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \leq \exp\left(n\mathbb{H}(\hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a} | \hat{\mathbf{U}}_{\mathcal{A}_b}, \hat{\mathbf{U}}_{\mathcal{B}_b})\right) \quad (8)$$

by [11, Lemma 2.5]. Letting  $K(\hat{\mathbf{U}}_{\mathcal{A}_b}, \hat{\mathbf{X}}_{\mathcal{A}_a}, \hat{\mathbf{X}}_{\mathcal{B}_a}, \hat{\mathbf{U}}_{\mathcal{B}_b})$  be the

number of elements of  $\mathfrak{S}$  with type  $(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b})$ ,

$$\begin{aligned} K(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) &= \sum_{(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \in \mathcal{F}} K(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \\ &\stackrel{(e)}{\leq} \sum_{(\mathbf{u}_{\mathcal{A}_b}, \mathbf{u}_{\mathcal{B}_b}) \in \mathcal{F}} \exp\left(n \mathbb{H}(\hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a} | \hat{U}_{\mathcal{A}_b} \hat{U}_{\mathcal{B}_b})\right) \\ &\stackrel{(f)}{\leq} \exp\left[n \left( \mathbb{I}(\mathbf{U}_{\mathcal{A}_b} \mathbf{U}_{\mathcal{B}_b}; \mathbf{X}_{\mathcal{A}_a} \mathbf{X}_{\mathcal{B}_a}) \right. \right. \\ &\quad \left. \left. + \mathbb{H}(\hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a} | \hat{U}_{\mathcal{A}_b} \hat{U}_{\mathcal{B}_b}) + \varepsilon \right)\right], \quad (9) \end{aligned}$$

where (e) follows from (8) and (f) from (5). Thus,

$$\begin{aligned} &\mathbb{P}\{(\bar{\mathbf{U}}_{\mathcal{A}_b}, \bar{\mathbf{U}}_{\mathcal{B}_b}) \in \mathcal{F}\} \\ &\stackrel{(g)}{\leq} \sum_{\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}} K(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) \\ &\quad \cdot \exp\left[-n \left( \mathbb{H}(\hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a}) + \text{D}_{\text{KL}}(\hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a} \| \bar{\mathbf{X}}_{\mathcal{A}_a} \bar{\mathbf{X}}_{\mathcal{B}_a}) \right)\right] \\ &\stackrel{(h)}{\leq} \sum_{\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}} \exp\left(-n(k(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) - \varepsilon)\right), \end{aligned}$$

where the sum is over all types that occur in  $\mathfrak{S}$ . Here, (g) follows from (7), and (h) from (9), and we defined

$$\begin{aligned} k(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) &\triangleq \mathbb{I}(\hat{U}_{\mathcal{A}_b} \hat{U}_{\mathcal{B}_b}; \hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a}) \\ &\quad - \mathbb{I}(\mathbf{U}_{\mathcal{A}_b} \mathbf{U}_{\mathcal{B}_b}; \mathbf{X}_{\mathcal{A}_a} \mathbf{X}_{\mathcal{B}_a}) + \text{D}_{\text{KL}}(\hat{X}_{\mathcal{A}_a} \hat{X}_{\mathcal{B}_a} \| \bar{\mathbf{X}}_{\mathcal{A}_a} \bar{\mathbf{X}}_{\mathcal{B}_a}). \end{aligned}$$

Using a type counting argument [11, Lemma 2.2], we can further bound the number of types that occur in  $\mathfrak{S}$  by

$$(n+1)^{|\mathcal{U}_{\mathcal{A}_b}| + |\mathcal{X}_{\mathcal{A}_a}| + |\mathcal{X}_{\mathcal{B}_a}| + |\mathcal{U}_{\mathcal{B}_b}|} \leq (n+1)^{|\mathcal{U}_{\mathcal{K}}| + |\mathcal{X}_{\mathcal{K}}|}$$

and obtain

$$\begin{aligned} &\mathbb{P}\{(\bar{\mathbf{U}}_{\mathcal{A}_b}, \bar{\mathbf{U}}_{\mathcal{B}_b}) \in \mathcal{F}\} \leq (n+1)^{|\mathcal{U}_{\mathcal{K}}| + |\mathcal{X}_{\mathcal{K}}|} \\ &\quad \times \max_{\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}} \exp\left[-n(k(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) - \varepsilon)\right], \quad (10) \end{aligned}$$

where the maximum is over all types  $(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b})$  occurring in  $\mathfrak{S}$ . For any such type, we have by construction  $(\hat{U}_{\mathcal{A}_b}, \hat{X}_{\mathcal{A}_a}, \hat{X}_{\mathcal{B}_a}, \hat{U}_{\mathcal{B}_b}) \in \mathcal{L}_{\delta}(\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b})$ . From (10), we can thus conclude:

$$\begin{aligned} &\mathbb{P}\{(\bar{\mathbf{U}}_{\mathcal{A}_b}, \bar{\mathbf{U}}_{\mathcal{B}_b}) \in \mathcal{F}\} \leq (n+1)^{|\mathcal{U}_{\mathcal{K}}| + |\mathcal{X}_{\mathcal{K}}|} \\ &\quad \times \max_{\mathcal{L}_{\delta}} \exp\left[-n(k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) - \varepsilon)\right], \quad (11) \end{aligned}$$

where we used  $\max_{\mathcal{L}_{\delta}}$  to denote maximization over all  $(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) \in \mathcal{L}_{\delta}(\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b})$ . Combining (6) and (11), we showed that for  $n$  large enough and some constant  $C$

$$\begin{aligned} \Theta(f_{\mathcal{A}}; f_{\mathcal{B}}) &\geq -\frac{\log(2)}{n} - \frac{1-\varepsilon}{n} \log \mathbb{P}\{(\bar{\mathbf{U}}_{\mathcal{A}_b}, \bar{\mathbf{U}}_{\mathcal{B}_b}) \in \mathcal{F}\} \\ &\geq -\varepsilon + (1-\varepsilon) \min_{\mathcal{L}_{\delta}} \left( k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) - \varepsilon \right) \\ &\geq -2\varepsilon + (1-\varepsilon) \min_{\mathcal{L}_{\delta}} \left( k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) \right) \\ &\geq \min_{\mathcal{L}_{\delta}} \left( k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) \right) - C\varepsilon, \quad (12) \end{aligned}$$

where  $\min_{\mathcal{L}_{\delta}}$  is used to denote minimization over all  $(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) \in \mathcal{L}_{\delta}(\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b})$ . Observing that  $k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b})$  is a continuous function of  $\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}$  and (12) holds for arbitrarily small  $\delta$ , we

obtain for  $n$  large enough,

$$\Theta(f_{\mathcal{A}}; f_{\mathcal{B}}) \geq \min_{\mathcal{L}_0} k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) - C'\varepsilon \quad (13)$$

for some (larger) constant  $C'$  by letting  $\delta \rightarrow 0$ . For  $(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) \in \mathcal{L}_0(\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b})$ , observe that we have

$$k(\tilde{U}_{\mathcal{A}_b}, \tilde{X}_{\mathcal{A}_a}, \tilde{X}_{\mathcal{B}_a}, \tilde{U}_{\mathcal{B}_b}) = \mathbb{I}(\tilde{U}_{\mathcal{A}_b} \tilde{X}_{\mathcal{A}_a}; \tilde{X}_{\mathcal{B}_a} \tilde{U}_{\mathcal{B}_b}). \quad (14)$$

and from the definition of  $\mathcal{L}_0(\mathbf{U}_{\mathcal{A}_b}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \mathbf{U}_{\mathcal{B}_b})$  clearly

$$\min_{\mathcal{L}_0} \mathbb{I}(\tilde{U}_{\mathcal{A}_b} \tilde{X}_{\mathcal{A}_a}; \tilde{X}_{\mathcal{B}_a} \tilde{U}_{\mathcal{B}_b}) \geq \mathbb{I}(\mathbf{U}_{\mathcal{A}_b}; \mathbf{U}_{\mathcal{B}_b}). \quad (15)$$

Combining (13), (14), and (15), we obtain  $\Theta(f_{\mathcal{A}}; f_{\mathcal{B}}) \geq \mathbb{I}(\mathbf{U}_{\mathcal{A}_b}; \mathbf{U}_{\mathcal{B}_b}) - C'\varepsilon \geq \mu_{\mathcal{A}, \mathcal{B}} - C'\varepsilon$  where the last inequality follows from (3). Thus,  $(\mu_{2\mathcal{A}, 2\mathcal{B}} - C'\varepsilon, R_{\mathcal{K}} + \varepsilon) \in \mathcal{R}$ , completing the proof since  $\varepsilon$  was arbitrary.

## V. SUMMARY AND DISCUSSION

We extended the information-theoretic biclustering problem to multiple sources and derived bounds on the achievable region. However, these bounds are not tight since the infamous Körner-Marton problem constitutes a counterexample. For the analog of the CEO problem we showed that our outer bound is tight in a special case, leveraging existing results from multi-terminal lossy source coding. The challenge of the biclustering problem lies in the fact that one needs to bound the mutual information between two arbitrary encodings solely based on their rates. This appears to be a difficult task since standard information-theoretic manipulations seem incapable of handling this dependence well.

## APPENDIX

### A. Proof of Lemma 11

Fix  $0 < \varepsilon', \varepsilon'' < \varepsilon$  and  $\tilde{R}_k = \mathbb{I}(\mathbf{X}_k; \mathbf{U}_k) + \varepsilon''/2$  for  $k \in \mathcal{K}$ .

**Encoding:** For  $n \in \mathbb{N}$  define  $\tilde{M}_k \triangleq e^{n\tilde{R}_k}$  and  $\tilde{\mathcal{M}}_k \triangleq \{1, 2, \dots, \tilde{M}_k\}$ . We apply [13, Lemma 3.4] and denote the random codebook  $\mathcal{C}_k \triangleq (\mathbf{V}_i^{(k)})_{i \in \tilde{\mathcal{M}}_k}$ , which are drawn independently uniform from  $\mathcal{T}_{[U_k]\delta}^{[n]}$  for each  $k \in \mathcal{K}$ . Denote the resulting randomized coding functions as  $\tilde{W}_k = \tilde{f}_k(\mathbf{X}_k, \mathcal{C}_k)$  and the corresponding decoded value as  $\tilde{\mathbf{U}}_k \triangleq \mathbf{V}_{\tilde{W}_k}^{(k)}$ . We have therefore, if  $n$  is chosen large enough and  $\delta$  small enough,

$$P_e \triangleq \mathbb{P}\left\{(\tilde{\mathbf{U}}_{\mathcal{K}}, \mathbf{X}_{\mathcal{K}}) \notin \mathcal{T}_{[U_{\mathcal{K}}]\delta}^{[n]}\right\} \leq \varepsilon'. \quad (16)$$

Next, we introduce (deterministic) binning. If  $R_k < \mathbb{I}(\mathbf{X}_k; \mathbf{U}_k)$ , partition  $\tilde{\mathcal{M}}_k$  into  $M_k \triangleq e^{n(R_k + \varepsilon'')}$  equally sized, consecutive bins, each of size  $e^{n\Delta_k}$  with  $\Delta_k \triangleq \tilde{R}_k - R_k - \varepsilon'' = \mathbb{I}(\mathbf{X}_k; \mathbf{U}_k) - R_k - \frac{\varepsilon''}{2}$ . The deterministic function  $\beta_k: \tilde{\mathcal{M}}_k \rightarrow \mathcal{M}_k$  maps a codeword onto the bin index in  $\mathcal{M}_k \triangleq \{1, 2, \dots, M_k\}$ , it belongs to. Now use the randomized encoding function  $f_k \triangleq \beta_k \circ \tilde{f}_k$ . If  $R_k \geq \mathbb{I}(\mathbf{X}_k; \mathbf{U}_k)$ , let  $\beta_k$  be the identity on  $\tilde{\mathcal{M}}_k$  and  $f_k \triangleq \tilde{f}_k$ .

**Decoding:** Given the codebooks, the decoding procedure  $g_{\mathcal{A}_a, \mathcal{A}_b}: \mathcal{M}_{\mathcal{A}_a} \rightarrow \mathcal{U}_{\mathcal{A}_b}^n$  for each  $\emptyset \neq \mathcal{A}_b \subseteq \mathcal{A}_a \subseteq \mathcal{K}$  is carried out as follows. Given  $m_{\mathcal{A}_a} \in \mathcal{M}_{\mathcal{A}_a}$ , let  $\tilde{m}_{\mathcal{A}_a} \triangleq \beta_{\mathcal{A}_a}^{-1}(m_{\mathcal{A}_a}) \subseteq \tilde{\mathcal{M}}_{\mathcal{A}_a}$  be all indices that are in the bins  $m_{\mathcal{A}_a}$ . Consider only the typical sequences  $\mathbf{V}_{\tilde{m}_{\mathcal{A}_a}}^{(\mathcal{A}_a)} \cap \mathcal{T}_{[U_{\mathcal{A}_a}]\delta}^n \triangleq \Phi \subseteq \mathcal{U}_{\mathcal{A}_a}^n$ . Denote the restriction of  $\Phi$  to the coordinates  $\mathcal{A}_b$  as  $[\Phi]_{\mathcal{A}_b}$ . If  $\Phi \neq \emptyset$ , choose the lexicographically smallest element of  $[\Phi]_{\mathcal{A}_b}$ , otherwise choose the lexicographically smallest element of  $\left[\mathbf{V}_{\tilde{m}_{\mathcal{A}_a}}^{(\mathcal{A}_a)}\right]_{\mathcal{A}_b}$ .

Let  $\mathcal{A}_a, \mathcal{A}_b, \mathcal{B}_a, \mathcal{B}_b \subset \mathcal{K}$  be sets of indices, such that the conditions of part 1 are satisfied. Using  $W_k \triangleq f_k(\mathbf{X}_k, \mathcal{C}_k)$  and the randomized decodings  $\tilde{\mathbf{U}}_1 \triangleq g_{\mathcal{A}_a, \mathcal{A}_b}(W_{\mathcal{A}_a}, \mathcal{C}_{\mathcal{A}_a})$  and  $\tilde{\mathbf{U}}_2 \triangleq g_{\mathcal{B}_a, \mathcal{B}_b}(W_{\mathcal{B}_a}, \mathcal{C}_{\mathcal{B}_a})$ , consider the error event  $\mathcal{E}_0 \triangleq \{(\tilde{\mathbf{U}}_1, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \tilde{\mathbf{U}}_2) \notin \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_b}, \mathcal{X}_{\mathcal{A}_a}, \mathcal{X}_{\mathcal{B}_a}, \mathcal{U}_{\mathcal{B}_b}]}^n\}$ . Define

$$\begin{aligned} \mathcal{E}_1 &\triangleq \{(\tilde{\mathbf{U}}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{A}_a}, \mathbf{X}_{\mathcal{B}_a}, \tilde{\mathbf{U}}_{\mathcal{B}_a}) \notin \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_b}, \mathcal{X}_{\mathcal{A}_a}, \mathcal{X}_{\mathcal{B}_a}, \mathcal{U}_{\mathcal{B}_b}]}^n\}, \\ \mathcal{E}_2 &\triangleq \left\{ \left| \left[ \mathbf{V}_{\tilde{\mathcal{W}}_{\mathcal{A}_a}}^{(\mathcal{A}_a)} \cap \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_a}]}^n \right]_{\mathcal{A}_b} \right| > 1 \right\}, \\ \mathcal{E}_3 &\triangleq \left\{ \left| \left[ \mathbf{V}_{\tilde{\mathcal{W}}_{\mathcal{B}_a}}^{(\mathcal{B}_a)} \cap \mathcal{T}_{[\mathcal{U}_{\mathcal{B}_a}]}^n \right]_{\mathcal{B}_b} \right| > 1 \right\}, \end{aligned}$$

with the random index set  $\tilde{\mathcal{W}}_{\mathcal{A}} \triangleq \beta_{\mathcal{A}}^{-1}(W_{\mathcal{A}})$ . We clearly have  $\mathcal{E}_0 \subseteq \mathcal{E}_1 \cup \mathcal{E}_2 \cup \mathcal{E}_3$  and thus

$$\begin{aligned} \mathbb{P}\{\mathcal{E}_0\} &\leq \mathbb{P}\{\mathcal{E}_1\} + \mathbb{P}\{\mathcal{E}_2|\mathcal{E}_1^c\} + \mathbb{P}\{\mathcal{E}_3|\mathcal{E}_1^c\} \\ &\stackrel{(j)}{\leq} \mathbb{P}\{\mathcal{E}_2|\mathcal{E}_1^c\} + \mathbb{P}\{\mathcal{E}_3|\mathcal{E}_1^c\} + \varepsilon', \end{aligned} \quad (17)$$

where (j) follows from (16). We can partition  $\tilde{\mathcal{W}}_{\mathcal{A}_a}$  into (random) subsets  $\mathcal{D}_{\mathcal{A}'}$ , indexed by  $\mathcal{A}' \subseteq \mathcal{A}_a$  as

$$\mathcal{D}_{\mathcal{A}'} \triangleq \{\tilde{w}_{\mathcal{A}_a} \in \tilde{\mathcal{W}}_{\mathcal{A}_a} : \tilde{w}_{\mathcal{A}'^c} = \tilde{W}_{\mathcal{A}'^c} \wedge \tilde{w}_k \neq \tilde{W}_k, \forall k \in \mathcal{A}'\},$$

where we used  $\mathcal{A}'^c \triangleq \mathcal{A}_a \setminus \mathcal{A}'$ . Observe that  $\mathcal{D}_{\emptyset} = \{\tilde{W}_{\mathcal{A}_a}\}$ . For each set  $\emptyset \neq \mathcal{A}' \subseteq \mathcal{A}_a$ , we define the error event

$$\mathcal{E}_{\mathcal{A}'} \triangleq \{\mathbf{V}_{\mathcal{D}_{\mathcal{A}'}}^{(\mathcal{A}_a)} \cap \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_a}]}^n \neq \emptyset\}$$

and we have

$$\mathbb{P}\{\mathcal{E}_2|\mathcal{E}_1^c\} \leq \sum_{\mathcal{A}' \subseteq \mathcal{A}_a : \mathcal{A}' \cap \mathcal{A}_b \neq \emptyset} \mathbb{P}\{\mathcal{E}_{\mathcal{A}'}|\mathcal{E}_1^c\} \quad (18)$$

since  $\mathcal{E}_2 \subseteq \bigcup_{\mathcal{A}' \subseteq \mathcal{A}_a : \mathcal{A}' \cap \mathcal{A}_b \neq \emptyset} \mathcal{E}_{\mathcal{A}'}$ . By the construction of the codebook,  $\mathcal{D}_{\mathcal{A}'}$  has  $\prod_{k \in \mathcal{A}'} (\exp[n\Delta_k] - 1)$  elements. For  $\tilde{w}_{\mathcal{A}_a} \in \mathcal{D}_{\mathcal{A}'}$ , the components  $\mathbf{V}_{\tilde{w}_{\mathcal{A}'}}^{(\mathcal{A}_a)}$  are uniformly distributed on  $\prod_{k \in \mathcal{A}'} \mathcal{T}_{[\mathcal{U}_k]}^n$  and  $\tilde{w}_{\mathcal{A}'^c} = \tilde{W}_{\mathcal{A}'^c}$ . Given  $\mathcal{E}_1^c$  we have in particular  $\tilde{\mathbf{U}}_{\mathcal{A}_a} \in \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_a}]}^n$ . Thus, for any  $\mathbf{u}_{\mathcal{A}'^c} \in \mathcal{T}_{[\mathcal{U}_{\mathcal{A}'^c}]}^n$ , we can conclude,

$$\begin{aligned} &\mathbb{P}\{\mathcal{E}_{\mathcal{A}'}|\mathcal{E}_1^c, \tilde{\mathbf{U}}_{\mathcal{A}'^c} = \mathbf{u}_{\mathcal{A}'^c}\} \\ &\leq \sum_{\tilde{w}_{\mathcal{A}_a} \in \mathcal{D}_{\mathcal{A}'}} \mathbb{P}\{\mathbf{V}_{\tilde{w}_{\mathcal{A}'}}^{(\mathcal{A}_a)} \in \mathcal{T}_{[\mathcal{U}_{\mathcal{A}_a}]}^n | \mathcal{E}_1^c, \tilde{\mathbf{U}}_{\mathcal{A}'^c} = \mathbf{u}_{\mathcal{A}'^c}\} \\ &\leq \exp \left[ n \left( \sum_{k \in \mathcal{A}'} \Delta_k \right) \right] \frac{|\mathcal{T}_{[\mathcal{U}_{\mathcal{A}'}|\mathcal{U}_{\mathcal{A}'^c}]}^n(\mathbf{u}_{\mathcal{A}'^c})|}{\prod_{k \in \mathcal{A}'} |\mathcal{T}_{[\mathcal{U}_k]}^n|} \\ &\stackrel{(k)}{\leq} \exp \left[ n \left( \varepsilon(\delta) + \mathbb{H}(\mathcal{U}_{\mathcal{A}'|\mathcal{U}_{\mathcal{A}'^c}}) + \sum_{k \in \mathcal{A}'} (\Delta_k - \mathbb{H}(\mathcal{U}_k)) \right) \right] \end{aligned} \quad (19)$$

where  $\varepsilon(\delta) = \sum_{k \in \mathcal{A}' \cup \emptyset} \varepsilon_k(\delta)$  goes to zero as  $\delta \rightarrow 0$  and (k) follows from the properties of types [10, Sections 2.4 and 2.5]. We observe that the definition of  $\tilde{R}_k$  and (1) imply for any  $\emptyset \neq \mathcal{A}' \subseteq \mathcal{A}_a$  with  $\mathcal{A}' \cap \mathcal{A}_b \neq \emptyset$ ,

$$\sum_{k \in \mathcal{A}'} \Delta_k \leq -\frac{\varepsilon''}{2} - \mathbb{H}(\mathcal{U}_{\mathcal{A}'|\mathcal{U}_{\mathcal{A}'^c}}) + \sum_{k \in \mathcal{A}'} \mathbb{H}(\mathcal{U}_k). \quad (20)$$

By taking the marginal distribution on  $\tilde{\mathbf{U}}_{\mathcal{A}'^c}$  in (19) and using (20), we obtain:

$$\mathbb{P}\{\mathcal{E}_{\mathcal{A}'}|\mathcal{E}_1^c\} \leq \exp \left[ n \left( \varepsilon(\delta) - \frac{\varepsilon''}{2} \right) \right] \leq \varepsilon', \quad (21)$$

for  $n$  large enough and  $\delta$  small enough. Similar reasoning for  $\mathbb{P}\{\mathcal{E}_3|\mathcal{E}_1^c\}$  and combination of (17), (18) and (21) yields

$$\mathbb{P}\{\mathcal{E}_0\} \leq \varepsilon' + 2^{|\mathcal{A}_a|} \varepsilon' + 2^{|\mathcal{B}_a|} \varepsilon' \leq 2^K \varepsilon'. \quad (22)$$

For a set  $\emptyset \neq \mathcal{A} \subseteq \mathcal{K}$ , we next analyze the random quantity  $\mathbb{L} \triangleq |\mathcal{C}_{\mathcal{A}} \cap \mathcal{T}_{[\mathcal{U}_{\mathcal{A}}]}^n|$ . For large  $n$  and  $\mathbf{V}_{\mathcal{A}} \in \mathcal{C}_{\mathcal{A}}$ , we have

$$\begin{aligned} \mathbb{E}[\mathbb{L}] &\leq \sum_{\mathbf{V}_{\mathcal{A}} \in \mathcal{C}_{\mathcal{A}}} \mathbb{E} \left[ \mathbb{1}_{\mathcal{T}_{[\mathcal{U}_{\mathcal{A}}]}^n}(\mathbf{V}_{\mathcal{A}}) \right] = \left( \prod_{k \in \mathcal{A}} \tilde{M}_k \right) \mathbb{E} \left[ \mathbb{1}_{\mathcal{T}_{[\mathcal{U}_{\mathcal{A}}]}^n}(\mathbf{V}_{\mathcal{A}}) \right] \\ &= \left( \prod_{k \in \mathcal{A}} \tilde{M}_k \right) \frac{|\mathcal{T}_{[\mathcal{U}_{\mathcal{A}}]}^n|}{\prod_{k \in \mathcal{A}} |\mathcal{T}_{[\mathcal{U}_k]}^n|} \\ &\stackrel{(l)}{\leq} \left( \prod_{k \in \mathcal{A}} \tilde{M}_k \right) \frac{\exp[n(\mathbb{H}(\mathcal{U}_{\mathcal{A}}) + \varepsilon_0(\delta))]}{\exp[n(\sum_{k \in \mathcal{A}} \mathbb{H}(\mathcal{U}_k) - \varepsilon_k(\delta))]} \\ &\leq \exp \left[ n \left( \mathbb{I}(\mathcal{U}_{\mathcal{A}}; \mathcal{X}_{\mathcal{A}}) + \hat{\varepsilon}(\delta) + |\mathcal{A}| \frac{\varepsilon''}{2} \right) \right], \end{aligned}$$

where  $\mathbb{1}_{\mathcal{Z}}(\cdot)$  denotes the indicator function of a set  $\mathcal{Z}$  and  $\hat{\varepsilon}(\delta) = \sum_{k \in \mathcal{A} \cup \emptyset} \varepsilon_k(\delta)$  vanishes to zero as  $\delta \rightarrow 0$ . Here, (l) follows from the properties of types [10, Sections 2.4 and 2.5]. We can choose  $\varepsilon''$  enough small s.t.  $\hat{\varepsilon}(\delta) + K\varepsilon''/2 < \varepsilon$  for suitably small  $\delta$ . Let  $\mathcal{E}_4 = \{\mathbb{L} \geq \exp[n(\mathbb{I}(\mathcal{U}_{\mathcal{A}}; \mathcal{X}_{\mathcal{A}}) + \varepsilon)]\}$  and from Markov's inequality provided that  $n$  is large enough:

$$\mathbb{P}\{\mathcal{E}_4\} \leq \exp \left[ n \left( \hat{\varepsilon}(\delta) - \varepsilon + |\mathcal{A}| \frac{\varepsilon''}{2} \right) \right] \leq \varepsilon'. \quad (23)$$

As  $\varepsilon'$  was arbitrary in (22) and (23), we can obtain deterministic encoding and decoding functions, s.t. (4) holds whenever the conditions of part 1 are satisfied. Taking into account that  $g_{\mathcal{A}_a, \mathcal{A}_b}(\mathcal{M}_{\mathcal{A}_a}) \times g_{\mathcal{B}_a, \mathcal{B}_b}(\mathcal{M}_{\mathcal{B}_a}) \subseteq \mathcal{C}_{\mathcal{A}_b \cup \mathcal{B}_b}$ , we also have (5). Notice that, given a specific code,  $\mathbb{P}\{\mathcal{E}_4\} < 1$  already implies  $\mathbb{P}\{\mathcal{E}_4\} = 0$  as the event  $\mathcal{E}_4$  is fully determined by the code  $\mathcal{C}_{\mathcal{K}}$  alone. ■

## REFERENCES

- [1] G. Pichler, P. Piantanida, and G. Matz, "Distributed information-theoretic biclustering of two memoryless sources," in *53rd Annual Allerton Conference on Communication, Control, and Computing*, 2015.
- [2] M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein, "Cluster analysis and display of genome-wide expression patterns," *Proc. of the Nat. Academy of Sciences*, vol. 95, no. 25, pp. 14 863–14 868, 1998.
- [3] K.-P. Wong, D. Feng, S. R. Meikle, and M. J. Fulham, "Segmentation of dynamic PET images using cluster analysis," *IEEE Trans. Nucl. Sci.*, vol. 49, no. 1, pp. 200–207, 2002.
- [4] G. Punj and D. W. Stewart, "Cluster analysis in marketing research: review and suggestions for application," *J. of Marketing Research*, vol. 20, no. 2, pp. 134–148, 1983.
- [5] T. S. Han, "Hypothesis testing with multiterminal data compression," *IEEE Trans. Inf. Theory*, vol. 33, no. 6, pp. 759–772, 1987.
- [6] T. S. Han and S. Amari, "Statistical inference under multiterminal data compression," *IEEE Trans. Inf. Theory*, vol. 44, no. 6, pp. 2300–2324, 1998.
- [7] J. Körner and K. Marton, "How to encode the modulo-two sum of binary sources," *IEEE Trans. Inf. Theory*, vol. 25, no. 2, pp. 219–221, March 1979.
- [8] T. A. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 740–761, Jan 2014.
- [9] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Wiley-Interscience, 2006.
- [10] A. El Gamal and Y.-H. Kim, *Network Information Theory*. Cambridge University Press, 2011.
- [11] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, 2011.
- [12] T. Berger, Z. Zhang, and H. Viswanathan, "The CEO problem," *IEEE Trans. Inf. Theory*, vol. 42, no. 3, pp. 887–902, May 1996.
- [13] T. S. Han and K. Kobayashi, "A unified achievable rate region for a general class of multiterminal source coding systems," *IEEE Trans. Inf. Theory*, vol. 26, no. 3, pp. 277–288, May 1980.