



# Information Bottleneck on General Alphabets

Georg Pichler, Günther Koliander

## ► To cite this version:

Georg Pichler, Günther Koliander. Information Bottleneck on General Alphabets. 2018 IEEE International Symposium on Information Theory (ISIT), Jun 2018, Vail, United States. <10.1109/ISIT.2018.8437714>. <hal-01930926>

**HAL Id: hal-01930926**

**<https://hal.science/hal-01930926v1>**

Submitted on 22 Nov 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# Information Bottleneck on General Alphabets

Georg Pichler, Günther Koliander

**Abstract**—We prove rigorously a source coding theorem that can probably be considered folklore, a generalization to arbitrary alphabets of a problem motivated by the Information Bottleneck method. For general random variables  $(Y, X)$ , we show essentially that for some  $n \in \mathbb{N}$ , a function  $f$  with rate limit  $\log|f| \leq nR$  and  $I(Y^n; f(X^n)) \geq nS$  exists if and only if there is a random variable  $U$  such that the Markov chain  $Y \dashv X \dashv U$  holds,  $I(U; X) \leq R$  and  $I(U; Y) \geq S$ . The proof relies on the well established discrete case and showcases a technique for lifting discrete coding theorems to arbitrary alphabets.

## I. INTRODUCTION

Since its inception [1], the *Information Bottleneck* (IB) method became a widely applied tool, especially in the context of machine learning problems. It has been successfully applied to various problems in machine learning [2], computer vision [3], and communications [4], [5], [6]. Furthermore, it is a valuable tool for channel output compression in a communication system [7], [8].

In the underlying information-theoretic problem, we define a pair  $(S, R) \in \mathbb{R}^2$  to be *achievable* for the two arbitrary random sources  $(Y, X)$ , if there exists a function  $f$  with rate limited range  $\frac{1}{n} \log|f| \leq R$  and  $I(Y; f(X)) \geq nS$ , where  $(Y, X)$  are  $n$  independent and identically distributed (i.i.d.) copies of  $(Y, X)$ .

While this Shannon-theoretic problem and variants thereof were also considered (e.g., [9], [10]), a large part of the literature is aimed at studying the IB function

$$S_{\text{IB}}(R) = \sup_{\substack{U : I(U; X) \leq R \\ Y \dashv X \dashv U}} I(U; Y) \quad (1)$$

in different contexts. In particular, several works (e.g., [1], [2], [11], [12], [13]) intend to compute a probability distribution that achieves the supremum in (1). The resulting distribution is then used as a building block in numerical algorithms, e.g., for document clustering [2] or dimensionality reduction [11].

In the discrete case,  $S_{\text{IB}}(R)$  is equal to the maximum of all  $S$  such that  $(S, R)$  is in the *achievable region* (closure of the set of all achievable pairs). This statement has been re-proven many times in different contexts [14], [10], [15], [16]. In this note, we prove a theorem, which can probably be considered folklore, extending this result from discrete to arbitrary random variables. Formally speaking, using the definitions in [17], we prove that a pair  $(S, R)$  is in the achievable region of an arbitrary source  $(Y, X)$  if and only if, for every  $\varepsilon > 0$ , there exists a random variable  $U$  with  $Y \dashv X \dashv U$ ,  $I(X; U) \leq R + \varepsilon$ , and  $I(Y; U) \geq S - \varepsilon$ . This provides a single-letter solution to the information-theoretic problem behind the information bottleneck method for arbitrary random sources and in particular it shows, that the information bottleneck for Gaussian random variables [11] is indeed the solution to a Shannon-theoretic problem.

The proof relies on the discrete case. Thus, the techniques employed could be useful for lifting other discrete coding theorems to the case of arbitrary alphabets.

## II. MAIN RESULT

Let  $Y$  and  $X$  be random variables with arbitrary alphabets  $\mathcal{S}_Y$  and  $\mathcal{S}_X$ , respectively. The bold-faced random vectors  $\mathbf{Y}$  and  $\mathbf{X}$  are  $n$  i.i.d. copies of  $Y$  and  $X$ , respectively. We then have the following definitions.

**Definition 1.** A pair  $(S, R) \in \mathbb{R}^2$  is *achievable* if for some  $n \in \mathbb{N}$  there exists a measurable function  $f: \mathcal{S}_X^n \rightarrow \mathcal{M}$  for some finite set  $\mathcal{M}$  with bounded cardinality  $\frac{1}{n} \log|\mathcal{M}| \leq R$  and

$$\frac{1}{n} I(\mathbf{Y}; f(\mathbf{X})) \geq S. \quad (2)$$

The set of all achievable pairs is denoted  $\mathcal{R} \subseteq \mathbb{R}^2$ .

**Definition 2.** A pair  $(S, R) \in \mathbb{R}^2$  is *IB-achievable* if there exists an additional random variable  $U$  with arbitrary alphabet  $\mathcal{S}_U$ , satisfying  $Y \dashv X \dashv U$  and

$$R \geq I(X; U), \quad (3)$$

$$S \leq I(Y; U). \quad (4)$$

The set of all IB-achievable pairs is denoted  $\mathcal{R}_{\text{IB}} \subseteq \mathbb{R}^2$ .

In what follows, we will prove the following theorem.

**Theorem 3.** The equality  $\overline{\mathcal{R}_{\text{IB}}} = \overline{\mathcal{R}}$  holds.

## III. PRELIMINARIES

When introducing a function, we implicitly assume it to be measurable w.r.t. the appropriate  $\sigma$ -algebras. The  $\sigma$ -algebra associated with a finite set is its power set and the  $\sigma$ -algebra associated with  $\mathbb{R}$  is the Borel  $\sigma$ -algebra. The symbol  $\emptyset$  is used for the empty set and for a constant random variable. When there is no possibility for confusion, we will not distinguish between a single-element set and its element, e.g., we write  $x$  instead of  $\{x\}$  and  $\mathbb{1}_x$  for the indicator function of  $\{x\}$ . We use  $A \triangle B := (A \setminus B) \cup (B \setminus A)$  to denote the symmetric set difference.

Let  $(\Omega, \Sigma, \mu)$  be a probability space. A random variable  $X: \Omega \rightarrow \mathcal{S}_X$  takes values in the measurable space  $(\mathcal{S}_X, \mathcal{A}_X)$ . The push-forward probability measure  $\mu_X: \mathcal{A}_X \rightarrow [0, 1]$  is defined by  $\mu_X(A) = \mu(X^{-1}(A))$  for all  $A \in \mathcal{A}_X$ . We will state most results in terms of push-forward measures and usually ignore the background probability space. When multiple random variables are defined, we implicitly assume the push-forward measures to be consistent in the sense that, e.g.,  $\mu_X(A) = \mu_{XY}(A \times \mathcal{S}_Y)$  for all  $A \in \mathcal{A}_X$ .

For  $n \in \mathbb{N}$  let  $\Omega^n$  denote the  $n$ -fold Cartesian product of  $(\Omega, \Sigma, \mu)$ . A bold-faced random vector, e.g.,  $\mathbf{X}$ , defined on  $\Omega^n$ , is an  $n$ -fold copy of  $X$ , i.e.,  $\mathbf{X} = X^n$ . Accordingly, the corresponding push-forward measure, e.g.,  $\mu_{\mathbf{X}}$  is the  $n$ -fold product measure.

For a random variable  $X$  let  $a_X$ ,  $b_X$ , and  $c_X$  denote arbitrary functions on  $\mathcal{S}_X$ , each with finite range. We will use the symbol  $\mathcal{M}_X$  to denote the range of  $a_X$ , i.e.,  $a_X: \mathcal{S}_X \rightarrow \mathcal{M}_X$ .

G. Pichler is with the Institute of Telecommunications, Technische Universität Wien, Vienna, Austria

G. Koliander is with the Acoustics Research Institute, Austrian Academy of Sciences, Vienna, Austria

Funding by WWTF Grants MA16-053, ICT15-119, and NXT17-013.

**Definition 4** ([18, Def. 8.11]). *The conditional expectation of a random variable  $X$  with  $S_X = \mathbb{R}$ , given a random variable  $Y$ , is a random variable  $\mathbb{E}[X|Y]$  such that*

- 1)  $\mathbb{E}[X|Y]$  is  $\sigma(Y)$ -measurable, and
- 2) for all  $A \in \sigma(Y)$ , we have  $\mathbb{E}[\mathbb{1}_A \mathbb{E}[X|Y]] = \mathbb{E}[\mathbb{1}_A X]$ .

The conditional probability of an event  $B \in \Sigma$  given  $Y$  is defined as  $P\{B|Y\} := \mathbb{E}[\mathbb{1}_B|Y]$ .

The conditional expectation and therefore also the conditional probability exists and is unique up to equality almost surely by [18, Thm. 8.12]. Furthermore, if  $(S_X, \mathcal{A}_X)$  is a standard space [17, Sec. 1.5], there even exists a *regular conditional distribution* of  $X$  given  $Y$  [18, Thm. 8.37].

**Definition 5.** *For two random variables  $X$  and  $Y$  a regular conditional distribution of  $X$  given  $Y$  is a function  $\kappa_{X|Y}: \Omega \times \mathcal{A}_X \rightarrow [0, 1]$  such that*

- 1) for every  $\omega \in \Omega$ , the set function  $\kappa_{X|Y}(\omega) := \kappa_{X|Y}(\omega; \cdot)$  is a probability measure on  $(S_X, \mathcal{A}_X)$ .
- 2) for every set  $A \in \mathcal{A}_X$ , the function  $\kappa_{X|Y}(\cdot; A)$  is  $\sigma(Y)$ -measurable.
- 3) for  $\mu$ -a. e.  $\omega \in \Omega$  and all  $A \in \mathcal{A}_X$ , we have  $\kappa_{X|Y}(\omega; A) = P\{X^{-1}(A)|Y\}(\omega)$  (cf. Def. 4).

Note, in particular, that finite spaces are standard spaces.

*Remark 1.* If the random variable  $Y$  is discrete, then  $\kappa_{X|Y}$  reduces to conditioning given events  $Y = y$  for  $y \in S_Y$ , i. e.,  $\kappa_{X|Y}(\omega; A) = \frac{\mu_{XY}(A \times Y(\omega))}{\mu_Y(Y(\omega))}$  (cf. [18, Lem. 8.10]).

We use the following definitions and results from [17], [18].

**Definition 6.** *For random variables  $X$  and  $Y$  with  $|S_X| < \infty$  the conditional entropy is defined as [17, Sec. 5.5]*

$$H(X|Y) := \int H(\kappa_{X|Y}) d\mu, \quad (5)$$

where  $H(\cdot)$  denotes discrete entropy on  $S_X$ . For arbitrary random variables  $X$ ,  $Y$ , and  $Z$  the conditional mutual information is defined as [17, Lem. 5.5.7]

$$\begin{aligned} I(X; Y|Z) &:= \sup_{a_X, a_Y} \int D(\kappa_{a_X(X)a_Y(Y)|Z} \| \kappa_{a_X(X)|Z} \times \kappa_{a_Y(Y)|Z}) d\mu \\ &= \sup_{a_X, a_Y} [H(a_X(X)|Z) + H(a_Y(Y)|Z) - H(a_X(X)a_Y(Y)|Z)], \end{aligned} \quad (6)$$

where  $D(\cdot \| \cdot)$  denotes Kullback-Leibler divergence [17, Sec. 2.3] and the supremum is taken over all  $a_X$  and  $a_Y$  with finite range. The mutual information is given by [17, Lem. 5.5.1]  $I(X; Y) := I(X; Y|\emptyset)$ .

**Definition 7** ([18, Def. 12.20]). *For arbitrary random variables  $X$ ,  $Y$ , and  $Z$ , the Markov chain  $X \dashv Y \dashv Z$  holds if, for any  $A \in \mathcal{A}_X$ ,  $B \in \mathcal{A}_Z$ , the following holds  $\mu$ -a. e.:*

$$P\{X^{-1}(A) \cap Z^{-1}(B)|Y\} = P\{X^{-1}(A)|Y\}P\{Z^{-1}(B)|Y\}. \quad (8)$$

In the following, we collect some properties of these definitions.

**Lemma 8.** *For random variables  $X$ ,  $Y$ , and  $Z$  the following properties hold:*

- (i)  $I(X; Y|Z) \geq 0$  with equality if and only if  $X \dashv Z \dashv Y$ .
- (ii) For discrete  $X$ , i. e.,  $|S_X| < \infty$ , we have  $I(X; Y) = H(X) - H(X|Y)$ .
- (iii)  $I(X; YZ) = I(X; Z) + I(X; Y|Z)$ .
- (iv) If  $X \dashv Y \dashv Z$ , then  $I(X; Y) \geq I(X; Z)$ .

*Proof.* (i): The claim  $I(X; Y|Z) \geq 0$  follows directly from (6) and the non-negativity of divergence.

Assume that  $X \dashv Z \dashv Y$ , i. e.,  $P\{X^{-1}(A) \cap Y^{-1}(B)|Z\} = P\{X^{-1}(A)|Z\}P\{Y^{-1}(B)|Z\}$  almost everywhere. Let  $a_X: S_X \rightarrow \mathcal{M}_X$  and  $a_Y: S_Y \rightarrow \mathcal{M}_Y$  be functions with finite range. Pick two arbitrary sets  $A \subseteq \mathcal{M}_X$ ,  $B \subseteq \mathcal{M}_Y$  and we obtain  $\mu$ -a. e.

$$\begin{aligned} \kappa_{a_X(X)a_Y(Y)|Z}(\cdot; A \times B) &= P\{X^{-1}(a_X^{-1}(A)) \cap Y^{-1}(a_Y^{-1}(B))|Z\} \end{aligned} \quad (9)$$

$$= P\{X^{-1}(a_X^{-1}(A))|Z\}P\{Y^{-1}(a_Y^{-1}(B))|Z\} \quad (10)$$

$$= \kappa_{a_X(X)|Z}(\cdot; A)\kappa_{a_Y(Y)|Z}(\cdot; B), \quad (11)$$

where (9) and (11) follow from part 3 of Def. 5. This proves that  $\mu$ -a. e. the equality of measures  $\kappa_{a_X(X)a_Y(Y)|Z} = \kappa_{a_X(X)|Z} \times \kappa_{a_Y(Y)|Z}$  holds. By the properties of Kullback-Leibler divergence [17, Thm. 2.3.1] we have  $I(X; Y|Z) = 0$  due to (6).

On the other hand, assume  $I(X; Y|Z) = 0$  and choose arbitrary sets  $A \in \mathcal{A}_X$  and  $B \in \mathcal{A}_Y$ . We define  $a_X := \mathbb{1}_A$ ,  $a_Y := \mathbb{1}_B$ ,  $\hat{X} := a_X(X)$ , and  $\hat{Y} := a_Y(Y)$ . By (6) we have  $D(\kappa_{\hat{X}\hat{Y}|Z}(\omega) \| \kappa_{\hat{X}|Z}(\omega) \times \kappa_{\hat{Y}|Z}(\omega)) = 0$  for  $\mu$ -a. e.  $\omega \in \Omega$ , which is equivalent to the equality  $\mu$ -a. e. of the measures  $\kappa_{\hat{X}\hat{Y}|Z} = \kappa_{\hat{X}|Z} \times \kappa_{\hat{Y}|Z}$ . We obtain  $\mu$ -a. e.,

$$P\{X^{-1}(A) \cap Y^{-1}(B)|Z\} = \kappa_{\hat{X}\hat{Y}|Z}(\cdot; 1 \times 1) \quad (12)$$

$$= \kappa_{\hat{X}|Z}(\cdot; 1)\kappa_{\hat{Y}|Z}(\cdot; 1) \quad (13)$$

$$= P\{X^{-1}(A)|Z\}P\{Y^{-1}(B)|Z\}. \quad (14)$$

(ii): See [17, Lem. 5.5.6].

(iii): See [17, Lem. 5.5.7].

(iv): Using Prop. (i) we have  $I(X; Z|Y) = 0$  and by Prop. (iii) it follows that

$$\begin{aligned} I(X; Z) &\leq I(X; YZ) \\ &= I(X; Y) + I(X; Z|Y) = I(X; Y). \end{aligned} \quad (15)$$

Occasionally we will interpret a probability measure on a finite space  $\mathcal{M}$  as a vector in  $[0, 1]^{\mathcal{M}}$ , equipped with the Borel  $\sigma$ -algebra. We will use the  $L_\infty$ -distance on this space.

**Definition 9.** *For two probability measures  $\mu$  and  $\nu$  on a finite space  $\mathcal{M}$ , their distance is defined as the  $L_\infty$ -distance  $d(\mu, \nu) := \max_{m \in \mathcal{M}} |\mu(m) - \nu(m)|$ . The diameter of  $A \subseteq [0, 1]^{\mathcal{M}}$  is defined as  $\text{diam}(A) = \sup_{\mu, \nu \in A} d(\mu, \nu)$ .*

**Lemma 10** ([19, Lem. 2.7]). *For two probability measures  $\mu$  and  $\nu$  on a finite space  $\mathcal{M}$  with  $d(\mu, \nu) \leq \varepsilon \leq \frac{1}{2}$  the inequality  $|H(\mu) - H(\nu)| \leq -\varepsilon |\mathcal{M}| \log \varepsilon$  holds.*

#### IV. PROOF OF $\mathcal{R}_{IB} \subseteq \overline{\mathcal{R}}$

For finite spaces  $S_Y$ ,  $S_X$ , and  $S_U$ , the statement  $\mathcal{R}_{IB} \subseteq \overline{\mathcal{R}}$  is well known, cf., [9, Sec. IV], [10, Sec. III.F]. We restate it in the form of the following lemma.

**Lemma 11.** *For random variables  $Y$ ,  $X$ , and  $U$  with finite  $S_Y$ ,  $S_X$ , and  $S_U$ , assume that  $Y \dashv X \dashv U$  holds. Then, for any  $\varepsilon > 0$ , there exists  $n \in \mathbb{N}$  and a function  $f: S_X^n \rightarrow \mathcal{M}$  with  $\frac{1}{n} \log |\mathcal{M}| \leq I(X; U) + \varepsilon$  such that  $\frac{1}{n} I(Y; f(X)) \geq I(Y; U) - \varepsilon$ .*

In a first step, we will utilize Lem. 11 to show  $\mathcal{R}_{IB} \subseteq \overline{\mathcal{R}}$  for an arbitrary alphabet  $S_X$ , i. e., we wish to prove the following Proposition 12, lifting the restriction  $|S_X| < \infty$ .

**Proposition 12.** *For random variables  $Y$ ,  $X$ , and  $U$  with finite  $S_Y$  and  $S_U$ , assume that  $Y \dashv X \dashv U$  holds. Then, for any*

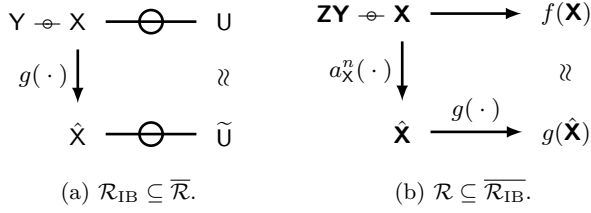


Fig. 1: Illustrations.

$\varepsilon > 0$ , there exists  $n \in \mathbb{N}$  and a function  $f: S_X^n \rightarrow \mathcal{M}$  with  $\frac{1}{n} \log |\mathcal{M}| \leq I(X; U) + \varepsilon$  such that

$$\frac{1}{n} I(\mathbf{Y}; f(\mathbf{X})) \geq I(Y; U) - \varepsilon. \quad (16)$$

*Remark 2.* Considering that both definitions of achievability (Defs. 1 and 2) only rely on the notion of mutual information, one may assume that Def. 6 can be used to directly infer Proposition 12 from Lem. 11. However, this is not the case. For an arbitrary discretization  $a_X(X)$  of  $X$ , we do have  $I(a_X(X); U) \leq I(X; U)$ . However, the Markov chain  $Y \circlearrowleft a_X(X) \circlearrowleft U$  does not hold in general. To circumvent this problem, we will use a discrete random variable  $\hat{X} = g(X)$  with an appropriate quantizer  $g$  and construct a new random variable  $\tilde{U}$ , satisfying the Markov chain  $Y \circlearrowleft \hat{X} \circlearrowleft \tilde{U}$  such that  $I(Y; \tilde{U})$  is close to  $I(Y; U)$ . Fig. 1a illustrates this strategy. We choose the quantizer  $g$  based on the conditional probability distribution of  $U$  given  $X$ , i. e., quantization based on  $\kappa_{U|X}$  using  $L_\infty$ -distance (cf. Def. 9). Subsequently, we will use that, by Lem. 10, a small  $L_\infty$ -distance guarantees a small gap in terms of information measures.

*Proof of Proposition 12.* Let  $\mu_{YXU}$  be a probability measure on  $\Omega := S_Y \times S_X \times S_U$ , such that  $Y \circlearrowleft X \circlearrowleft U$  holds. Fix  $0 < \delta \leq \frac{1}{2}$  and find a finite, measurable partition  $(P_i)_{i \in \mathcal{I}}$  of the space of probability measures on  $S_U$  such that for every  $i \in \mathcal{I}$  we have  $\text{diam}(P_i) \leq \delta$  and fix some  $\nu_i \in P_i$  for every  $i \in \mathcal{I}$ . Define the random variable  $\hat{X}: \Omega \rightarrow \mathcal{I}$  as  $\hat{X} = i$  if  $\kappa_{U|X} \in P_i$ . The random variable  $\hat{X}$  is  $\sigma(X)$ -measurable (see Appendix A). We can therefore find a measurable function  $g$  such that  $\hat{X} = g(X)$  by the factorization lemma [18, Corollary 1.97]. Define the new probability space  $\Omega \times \times_{i \in \mathcal{I}} S_U$ , equipped with the probability measure  $\mu_{YXU\tilde{U}} := \mu_{YXU} \times \times_{i \in \mathcal{I}} \nu_i$ . Slightly abusing notation, we define the random variables  $Y, X, U$ , and  $\tilde{U}_i$  (for every  $i \in \mathcal{I}$ ) as the according projections. We also use  $\hat{X} = g(X)$  and define the random variable  $\tilde{U} = \tilde{U}_{\hat{X}}$ . From this construction we have  $\mu_{YXU\tilde{U}}$ -a. e. the equality of measures  $\kappa_{\tilde{U}|\hat{X}} = \kappa_{\tilde{U}|X} = \nu_{\hat{X}}$ , as well as  $Y \circlearrowleft \hat{X} \circlearrowleft \tilde{U}$  and  $Y \circlearrowleft X \circlearrowleft \tilde{U}$  (see Appendix B). Therefore, we have  $\mu_{YXU\tilde{U}}$ -a. e.

$$d(\kappa_{\tilde{U}|\hat{X}}, \kappa_{U|X}) \leq \delta, \text{ and } d(\kappa_{\tilde{U}|\hat{X}}, \kappa_{U|X}) \leq \delta, \quad (17)$$

by  $\kappa_{\tilde{U}|\hat{X}} = \kappa_{\tilde{U}|X} = \nu_{\hat{X}}$  and  $\kappa_{U|X}, \nu_{\hat{X}} \in P_{\hat{X}}$ . Thus, for any  $u \in S_U$ ,

$$\mu_U(u) = \int \kappa_{U|X}(\cdot; u) d\mu_{YXU} \quad (18)$$

$$\leq \int (\kappa_{\tilde{U}|\hat{X}}(\cdot; u) + \delta) d\mu_{YXU\tilde{U}} = \mu_{\tilde{U}}(u) + \delta \quad (19)$$

and, by the same argument,  $\mu_U(u) \geq \mu_{\tilde{U}}(u) - \delta$ , i. e., in total,

$$d(\mu_U, \mu_{\tilde{U}}) \leq \delta. \quad (20)$$

Thus, we obtain

$$I(X; U) = H(\mu_U) - H(U|X) \quad (21)$$

$$\stackrel{(20)}{\geq} H(\mu_{\tilde{U}}) + \delta |S_U| \log \delta - \int H(\kappa_{U|X}) d\mu_{YXU} \quad (22)$$

$$\stackrel{(17)}{\geq} H(\mu_{\tilde{U}}) + 2\delta |S_U| \log \delta - \int H(\kappa_{\tilde{U}|\hat{X}}) d\mu_{YXU\tilde{U}} \quad (23)$$

$$= I(\hat{X}; \tilde{U}) + 2\delta |S_U| \log \delta, \quad (24)$$

where (21) and (24) follow from Prop. (ii) of Lem. 8, and in both (22) and (23) we used Lem. 10. From  $Y \circlearrowleft X \circlearrowleft U$  and Prop. (i) of Lem. 8, we know that  $\mu_{YXU}$ -a. e., we have the equality of measures  $\kappa_{YU|X} = \kappa_{Y|X} \times \kappa_{U|X}$ . Using this equality in (26) we obtain

$$\mu_{YU}(y \times u) = \int \kappa_{YU|X}(\cdot; y \times u) d\mu_{YXU} \quad (25)$$

$$= \int \kappa_{Y|X}(\cdot; y) \kappa_{U|X}(\cdot; u) d\mu_{YXU} \quad (26)$$

$$\stackrel{(17)}{\leq} \int \kappa_{Y|X}(\cdot; y) (\kappa_{\tilde{U}|\hat{X}}(\cdot; u) + \delta) d\mu_{YXU\tilde{U}} \quad (27)$$

$$\leq \int \kappa_{\tilde{U}|\hat{X}}(\cdot; y \times u) d\mu_{YXU\tilde{U}} + \delta \quad (28)$$

$$= \mu_{\tilde{U}}(y \times u) + \delta, \quad (29)$$

where (25) and (29) follow from the defining property of conditional probability, part 2 of Def. 4, and (28) follows from  $Y \circlearrowleft X \circlearrowleft \tilde{U}$  and Prop. (i) of Lem. 8. By the same argument, one can show that  $\mu_{YU}(y \times u) \geq \mu_{\tilde{U}}(y \times u) - \delta$ . Therefore, in total,  $d(\mu_{YU}, \mu_{\tilde{U}}) \leq \delta$  and, by Lem. 10,

$$|H(YU) - H(Y\tilde{U})| \leq -\delta |S_Y| |S_U| \log \delta. \quad (30)$$

Thus, the mutual information can be bounded by

$$I(Y; U) = H(Y) + H(U) - H(YU) \quad (31)$$

$$\stackrel{(20)}{\leq} H(Y) + H(\tilde{U}) - \delta |S_U| \log \delta - H(YU) \quad (32)$$

$$\stackrel{(30)}{\leq} I(Y; \tilde{U}) - \delta (|S_Y| + 1) |S_U| \log \delta \quad (33)$$

$$\leq I(Y; \tilde{U}) - 2\delta |S_Y| |S_U| \log \delta, \quad (34)$$

where we applied Lem. 10 in (32) and (33). We apply Lem. 11 to the three random variables  $Y, \hat{X}$ , and  $\tilde{U}$  and obtain a function  $\hat{f}: \mathcal{I}^n \rightarrow \mathcal{M}$  with  $\frac{1}{n} I(\mathbf{Y}; \hat{f}(\hat{\mathbf{X}})) \geq I(Y; \tilde{U}) - \delta$  and

$$\frac{1}{n} \log |\mathcal{M}| \leq I(\hat{X}; \tilde{U}) + \delta \stackrel{(24)}{\leq} I(X; U) + \delta - 2\delta |S_U| \log \delta. \quad (35)$$

We have  $\hat{\mathbf{X}} = g^n \circ \mathbf{X}$  and defining  $f := \hat{f} \circ g^n$ , we obtain

$$\frac{1}{n} I(\mathbf{Y}; f(\mathbf{X})) = \frac{1}{n} I(\mathbf{Y}; \hat{f}(\hat{\mathbf{X}})) \geq I(Y; \tilde{U}) - \delta \quad (36)$$

$$\stackrel{(34)}{\geq} I(Y; U) + 2\delta |S_Y| |S_U| \log \delta - \delta. \quad (37)$$

Choosing  $\delta$  such that  $\varepsilon \geq -2\delta |S_Y| |S_U| \log \delta + \delta$  completes the proof.  $\blacksquare$

We can now complete the proof by showing the following lemma.

**Lemma 13.**  $\mathcal{R}_{IB} \subseteq \overline{\mathcal{R}}$ .

*Proof.* Assuming  $(S, R) \in \mathcal{R}_{IB}$ , choose  $\mu_{YXU}$  according to Def. 2. Clearly  $I(X; U) < \infty$  to satisfy (3) and thus also  $I(Y; U) < \infty$  by Prop. (iv) of Lem. 8 as  $Y \circlearrowleft X \circlearrowleft U$  holds. Pick  $\varepsilon > 0$ , select

functions  $a_X, a_U$  such that  $I(a_X(X); a_U(U)) \geq I(X; U) - \varepsilon$ , and select functions  $b_Y, b_U$  such that  $I(b_Y(Y); b_U(U)) \geq I(Y; U) - \varepsilon$  (cf. (7)). Using  $\hat{U} := (a_U(U), b_U(U))$  and  $\hat{Y} := b_Y(Y)$ , we have

$$0 = I(Y; U|X) = \sup_{c_Y, c_U} I(c_Y(Y); c_U(U)|X) \geq I(\hat{Y}; \hat{U}|X) \geq 0 \quad (38)$$

as well as

$$I(X; U) = \sup_{c_X, c_U} I(c_X(X); c_U(U)) \quad (39)$$

$$\geq \sup_{c_X} I(c_X(X); \hat{U}) = I(X; \hat{U}), \text{ and} \quad (40)$$

$$I(Y; U) - \varepsilon \leq I(b_Y(Y); b_U(U)) \leq I(\hat{Y}; \hat{U}). \quad (41)$$

We apply Proposition 12, substituting  $\hat{U} \rightarrow U$  and  $\hat{Y} \rightarrow Y$ . Proposition 12 guarantees the existence of a function  $f: S_X^n \rightarrow \mathcal{M}$  with  $\frac{1}{n} \log |\mathcal{M}| \leq I(X; \hat{U}) + \varepsilon \leq I(X; U) + \varepsilon \leq R + \varepsilon$  and

$$\frac{1}{n} I(Y; f(\mathbf{X})) = \frac{1}{n} \sup_{c_Y} I(c_Y \circ Y; f(\mathbf{X})) \quad (42)$$

$$\geq \frac{1}{n} I(b_Y^n \circ Y; f(\mathbf{X})) = \frac{1}{n} I(\hat{Y}; f(\mathbf{X})) \quad (43)$$

$$\stackrel{(16)}{\geq} I(\hat{Y}; \hat{U}) - \varepsilon \stackrel{(41)}{\geq} I(Y; U) - 2\varepsilon \stackrel{(4)}{\geq} S - 2\varepsilon. \quad (44)$$

Thus,  $(S - 2\varepsilon, R - \varepsilon) \in \mathcal{R}$  and therefore  $(S, R) \in \overline{\mathcal{R}}$ . ■

#### V. PROOF OF $\mathcal{R} \subseteq \overline{\mathcal{R}}_{IB}$

We start with the well-known result  $\mathcal{R}_{IB} \subseteq \overline{\mathcal{R}}$  for finite spaces  $S_Y, S_X$ , and  $S_U$ , cf., [9, Sec. IV], [10, Sec. III.F]. The statement is rephrased in the following lemma.

**Lemma 14.** *Assume that the spaces  $S_Y$  and  $S_X$  are both finite and  $\mu_{YX}$  is fixed. For some  $n \in \mathbb{N}$ , let  $f: S_X^n \rightarrow \mathcal{M}$  be a function with  $|\mathcal{M}| < \infty$ . Then there exists a probability measure  $\mu_{YXU}$ , extending  $\mu_{YX}$ , such that  $S_U$  is finite,  $Y \circlearrowleft X \circlearrowleft U$ , and*

$$I(X; U) \leq \frac{1}{n} \log |\mathcal{M}|, \quad (45)$$

$$I(Y; U) \geq \frac{1}{n} I(Y; f(\mathbf{X})). \quad (46)$$

We can slightly strengthen Lem. 14.

**Corollary 15.** *Assume that, in the setting of Lem. 14, we are given  $\mu_{ZYX}$  on  $S_Z \times S_Y \times S_X$ , extending  $\mu_{YX}$ , where  $S_Z$  is arbitrary, not necessarily finite. Then there exists a probability measure  $\mu_{ZYXU}$ , extending  $\mu_{ZYX}$ , such that  $S_U$  is finite and  $ZY \circlearrowleft X \circlearrowleft U$ , (45), and (46) hold.*

*Proof.* Apply Lem. 14 to obtain  $\mu_{YXU}$  on  $S_Y \times S_X \times S_U$  satisfying (45), (46), and  $Y \circlearrowleft X \circlearrowleft U$ . We define  $\mu_{ZYXU}$  by

$$\mu_{ZYXU}(A \times y \times x \times u) = \frac{\mu_{ZYX}(A \times y \times x)}{\mu_{YX}(y \times x)} \mu_{YXU}(y \times x \times u) \quad (47)$$

for any  $(y, x, u) \in S_Y \times S_X \times S_U$  and  $A \in \mathcal{A}_Z$ . Pick arbitrary  $A \in \mathcal{A}_Z$ ,  $y \in S_Y$ , and  $u \in S_U$ . The Markov chain  $ZY \circlearrowleft X \circlearrowleft U$  now follows as the events  $Z^{-1}(A) \cap Y^{-1}(y)$  and  $U^{-1}(u)$  are independent given  $X^{-1}(x)$  for any  $x \in S_X$  (cf. Rmk. 1). ■

Again, we proceed by extending Cor. 15, lifting the restriction that  $S_X$  is finite and obtain the following proposition.

**Proposition 16.** *Given a probability measure  $\mu_{ZYX}$  as in Cor. 15, assume that  $|\mathcal{S}_Y| < \infty$ . For some  $n \in \mathbb{N}$ , let  $f: S_X^n \rightarrow \mathcal{M}$  be a function with  $|\mathcal{M}| < \infty$ . Then, for any  $\varepsilon > 0$ , there exists a*

*probability measure  $\mu_{ZYXU}$ , extending  $\mu_{ZYX}$  with  $ZY \circlearrowleft X \circlearrowleft U$  and*

$$I(X; U) \leq \frac{1}{n} \log |\mathcal{M}| \quad (48)$$

$$I(Y; U) \geq \frac{1}{n} I(Y; f(\mathbf{X})) - \varepsilon. \quad (49)$$

*Remark 3.* In contrast to Proposition 12, Proposition 16 could be proved by the usual single-letterization + time-sharing strategy, by showing that the necessary Markov chains hold. However, we will rely on the discrete case (Lem. 14) and showcase a technique to lift it to general alphabets.

*Remark 4.* In the proof of Proposition 16, we face a similar problem as outlined in Rmk. 2. We need to construct a function  $g(\hat{\mathbf{X}})$  of a “per-letter” quantization  $\hat{\mathbf{X}} := a_X^n(\mathbf{X})$ , that is close to  $f(\mathbf{X})$  in distribution. Fig. 1b provides a sketch.

*Proof of Proposition 16.* We can partition  $S_X^n = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$  into finitely many measurable, mutually disjoint sets  $\mathcal{Q}_m := f^{-1}(m)$ ,  $m \in \mathcal{M}$ . We want to approximate the sets  $\mathcal{Q}_m$  by a finite union of rectangles in the semiring [18, Def. 1.9]  $\Xi := \{\mathcal{B} : \mathcal{B} = \bigtimes_{i=1}^n B_i \text{ with } B_i \in \mathcal{A}_X\}$ . We choose  $\delta > 0$ , which will be specified later. According to [18, Thm. 1.65(ii)], we obtain  $\mathcal{B}^{(m)} := \bigcup_{k=1}^K \mathcal{B}_k^{(m)}$  for each  $m \in \mathcal{M}$ , where  $\mathcal{B}_k^{(m)} \in \Xi$  are mutually disjoint sets, satisfying  $\mu_X(\mathcal{B}^{(m)} \Delta \mathcal{Q}_m) \leq \delta$ . Since  $\mathcal{B}_k^{(m)} \in \Xi$ , we have  $\mathcal{B}_k^{(m)} = \bigtimes_{i=1}^n B_{k,i}^{(m)}$  for some  $B_{k,i}^{(m)} \in \mathcal{A}_X$ . We can construct functions  $a_X$  and  $g$  such that  $g \circ a_X^n(\mathbf{x}) = m$  whenever  $\mathbf{x} \in \mathcal{B}^{(m)}$  and  $\mathbf{x} \notin \mathcal{B}^{(m')}$  with  $\mathcal{B}^{(m')} := \bigcup_{m' \neq m} \mathcal{B}^{(m')}$ . Indeed, we obtain  $a_X$  by finding a measurable partition of  $S_X$  that is finer than  $(B_{k,i}^{(m)}, (B_{k,i}^{(m)})^c)$  for all  $i, k, m$ . For fixed  $m \in \mathcal{M}$ ,

$$\mathcal{Q}_m \subseteq \mathcal{Q}_m \cup (\mathcal{B}^{(m)} \setminus \mathcal{B}^{(m')}) \quad (50)$$

$$\subseteq (\mathcal{B}^{(m)} \setminus \mathcal{B}^{(m')}) \cup (\mathcal{Q}_m \setminus \mathcal{B}^{(m)}) \cup \bigcup_{m' \neq m} \mathcal{Q}_m \cap \mathcal{B}^{(m')} \quad (51)$$

$$\subseteq (\mathcal{B}^{(m)} \setminus \mathcal{B}^{(m')}) \cup (\mathcal{Q}_m \Delta \mathcal{B}^{(m)}) \cup \bigcup_{m' \neq m} \mathcal{B}^{(m')} \setminus \mathcal{Q}_{m'} \quad (52)$$

$$\subseteq (\mathcal{B}^{(m)} \setminus \mathcal{B}^{(m')}) \cup \bigcup_{m'} \mathcal{B}^{(m')} \Delta \mathcal{Q}_{m'}, \quad (53)$$

where we used the fact that  $\mathcal{Q}_m \cap \mathcal{Q}_{m'} = \emptyset$  for  $m \neq m'$  in (52). Using  $\hat{X} := a_X(X)$ , we obtain for any  $y \in S_Y^n$

$$\mu_{Yf(\mathbf{X})}(y \times m) = \mu_{YX}(y \times \mathcal{Q}_m) \quad (54)$$

$$\stackrel{(53)}{\leq} \mu_{YX}(y \times (\mathcal{B}^{(m)} \setminus \mathcal{B}^{(m')})) + \sum_{m'} \mu_X(\mathcal{B}^{(m')} \Delta \mathcal{Q}_{m'}) \quad (55)$$

$$\leq \mu_{Yg(\hat{\mathbf{X}})}(y \times m) + |\mathcal{M}| \delta. \quad (56)$$

On the other hand, we have

$$\mu_{Yf(\mathbf{X})}(y \times m) = \mu_Y(y) - \sum_{m' \neq m} \mu_{Yf(\mathbf{X})}(y \times m') \quad (57)$$

$$\stackrel{(56)}{\geq} \mu_Y(y) - \sum_{m' \neq m} (\mu_{Yg(\hat{\mathbf{X}})}(y \times m') + |\mathcal{M}| \delta) \quad (58)$$

$$\geq \mu_{Yg(\hat{\mathbf{X}})}(y \times m) - |\mathcal{M}|^2 \delta. \quad (59)$$

We thus obtain  $d(\mu_{Yf(\mathbf{X})}, \mu_{Yg(\hat{\mathbf{X}})}) \leq |\mathcal{M}|^2 \delta$ . This also implies  $d(\mu_{f(\mathbf{X})}, \mu_{g(\hat{\mathbf{X}})}) \leq |\mathcal{S}_Y|^n |\mathcal{M}|^2 \delta$ . Assume  $|\mathcal{S}_Y|^n |\mathcal{M}|^2 \delta \leq \frac{1}{2}$  and apply Cor. 15 substituting  $\hat{X} \rightarrow X$ ,  $XZ \rightarrow Z$ , and the function  $g \rightarrow f$ . This yields a random variable  $U$  with  $XZY \circlearrowleft \hat{X} \circlearrowleft U$ ,

$$I(\hat{X}; U) \leq \frac{1}{n} \log |\mathcal{M}|, \text{ and } I(Y; U) \geq \frac{1}{n} I(Y; g(\hat{\mathbf{X}})). \quad (60)$$

We also obtain  $\mathbf{ZY} \dashv\!\!\!\dashv \mathbf{X} \dashv\!\!\!\dashv \mathbf{U}$  due to

$$0 = \mathbf{I}(\mathbf{XZY}; \mathbf{U}|\hat{\mathbf{X}}) \quad (61)$$

$$= \mathbf{I}(\mathbf{XZY}; \mathbf{U}) - \mathbf{I}(\mathbf{U}; \hat{\mathbf{X}}) \quad (62)$$

$$\geq \mathbf{I}(\mathbf{XZY}; \mathbf{U}) - \mathbf{I}(\mathbf{U}; \mathbf{X}) \quad (63)$$

$$= \mathbf{I}(\mathbf{ZY}; \mathbf{U}|\mathbf{X}) \quad (64)$$

$$\geq 0, \quad (65)$$

where (61) follows from  $\mathbf{XZY} \dashv\!\!\!\dashv \hat{\mathbf{X}} \dashv\!\!\!\dashv \mathbf{U}$  using Prop. (i) of Lem. 8, (62) and (64) follow from Prop. (iii) of Lem. 8, (63) is a consequence of Def. 6, and we used Prop. (i) of Lem. 8 in (65). This also immediately implies  $0 = \mathbf{I}(\mathbf{X}; \mathbf{U}|\hat{\mathbf{X}})$  and hence

$$\frac{1}{n} \log |\mathcal{M}| \stackrel{(60)}{\geq} \mathbf{I}(\hat{\mathbf{X}}; \mathbf{U}) = \mathbf{I}(\hat{\mathbf{X}}; \mathbf{U}) + \mathbf{I}(\mathbf{X}; \mathbf{U}|\hat{\mathbf{X}}) \quad (66)$$

$$= \mathbf{I}(\mathbf{X}\hat{\mathbf{X}}; \mathbf{U}) = \mathbf{I}(\mathbf{X}; \mathbf{U}), \quad (67)$$

where we used Prop. (iii) of Lem. 8 in (67). We also have

$$\mathbf{I}(\mathbf{Y}; \mathbf{U}) \stackrel{(60)}{\geq} \frac{1}{n} \mathbf{I}(\mathbf{Y}; g(\hat{\mathbf{X}})) \quad (68)$$

$$= \frac{1}{n} (\mathbf{H}(\mathbf{Y}) + \mathbf{H}(g(\hat{\mathbf{X}})) - \mathbf{H}(\mathbf{Y}g(\hat{\mathbf{X}}))) \quad (69)$$

$$\geq \frac{1}{n} \mathbf{I}(\mathbf{Y}; f(\mathbf{X})) + \frac{1}{n} |\mathcal{S}_Y|^n |\mathcal{M}|^3 \delta \log(|\mathcal{M}|^2 \delta) + \frac{1}{n} |\mathcal{S}_Y|^n |\mathcal{M}|^3 \delta \log(|\mathcal{S}_Y|^n |\mathcal{M}|^2 \delta) \quad (70)$$

$$\geq \frac{1}{n} \mathbf{I}(\mathbf{Y}; f(\mathbf{X})) + \frac{2}{n} |\mathcal{S}_Y|^n |\mathcal{M}|^3 \delta \log(|\mathcal{M}|^2 \delta) \quad (71)$$

where we used Lem. 10 in (70). Select  $\delta$  such that  $\varepsilon \geq -\frac{2}{n} |\mathcal{S}_Y|^n |\mathcal{M}|^3 \delta \log(|\mathcal{M}|^2 \delta)$ . ■

We can now finish the proof by showing the following lemma.

**Lemma 17.**  $\mathcal{R} \subseteq \overline{\mathcal{R}_{\text{IB}}}$ .

*Proof.* Assume  $(S, R) \in \mathcal{R}$  and choose  $n \in \mathbb{N}$  and  $f$ , satisfying  $\frac{1}{n} \log |\mathcal{M}| \leq R$  and (2). Choose any  $\varepsilon > 0$  and find  $a_Y$  such that

$$\mathbf{I}(a_Y^n(\mathbf{Y}); f(\mathbf{X})) \geq \mathbf{I}(\mathbf{Y}; f(\mathbf{X})) - \varepsilon \stackrel{(2)}{\geq} nS - \varepsilon. \quad (72)$$

This is possible by applying [17, Lem. 5.2.2] with the algebra that is generated by the rectangles (cf. the paragraph above [17, Lem. 5.5.1]). We apply Proposition 16, substituting  $a_Y(\mathbf{Y}) \rightarrow \mathbf{Y}$  and  $\mathbf{Y} \rightarrow \mathbf{Z}$ . For arbitrary  $\varepsilon > 0$ , Proposition 16 provides  $\mathbf{U}$  with  $\mathbf{Y}a_Y(\mathbf{Y}) \dashv\!\!\!\dashv \mathbf{X} \dashv\!\!\!\dashv \mathbf{U}$  (i.e.,  $\mathbf{Y} \dashv\!\!\!\dashv \mathbf{X} \dashv\!\!\!\dashv \mathbf{U}$ ) and

$$\mathbf{I}(\mathbf{X}; \mathbf{U}) \leq \frac{1}{n} \log |\mathcal{M}| \leq R \quad (73)$$

$$\mathbf{I}(\mathbf{Y}; \mathbf{U}) \geq \mathbf{I}(a_Y(\mathbf{Y}); \mathbf{U}) \quad (74)$$

$$\stackrel{(49)}{\geq} \frac{1}{n} \mathbf{I}(a_Y^n(\mathbf{Y}); f(\mathbf{X})) - \varepsilon \stackrel{(72)}{\geq} S - 2\varepsilon. \quad (75)$$

Hence,  $(S - 2\varepsilon, R) \in \mathcal{R}_{\text{IB}}$  and consequently  $(S, R) \in \overline{\mathcal{R}_{\text{IB}}}$ . ■

## APPENDIX

### A. $\hat{\mathbf{X}}$ is $\sigma(\mathbf{X})$ -measurable

For  $u \in \mathcal{S}_U$  consider the  $\sigma(\mathbf{X})$ -measurable function  $h_u := \kappa_{U|\mathbf{X}}(\cdot; u)$  on  $[0, 1]$ . We obtain the vector valued function  $h := (h_u)_{u \in \mathcal{S}_U}$  on  $[0, 1]^{|\mathcal{S}_U|}$ . This function  $h$  is  $\sigma(\mathbf{X})$ -measurable as every component is  $\sigma(\mathbf{X})$ -measurable. Thus, we have  $\hat{\mathbf{X}}^{-1}(i) = h^{-1}(P_i) \in \sigma(\mathbf{X})$ .

### B. Distribution of $\tilde{\mathbf{U}}$ and Conditional Independence

We will first show that  $\mu_{\mathbf{YXU}\tilde{\mathbf{U}}_T}$ -a.e.

$$\kappa_{\tilde{\mathbf{U}}|\hat{\mathbf{X}}} = \kappa_{\tilde{\mathbf{U}}|\mathbf{X}} = \nu_{\hat{\mathbf{X}}}. \quad (76)$$

Clearly,  $\nu_{\hat{\mathbf{X}}}$  is a probability measure everywhere. Fixing  $u \in \mathcal{S}_U$ , we need that  $\nu_{\hat{\mathbf{X}}}(u)$  is  $\sigma(\hat{\mathbf{X}})$ -measurable, which is shown by the factorization lemma [18, Corollary 1.97], when writing  $\nu_{\hat{\mathbf{X}}}(u) = \nu_{(\cdot)}(u) \circ \hat{\mathbf{X}}$ . Also, this proves  $\sigma(\mathbf{X})$ -measurability as  $\hat{\mathbf{X}}$  is  $\sigma(\mathbf{X})$ -measurable, i.e.,  $\sigma(\hat{\mathbf{X}}) \subseteq \sigma(\mathbf{X})$ . It remains to show the defining property of conditional probability, part 2 of Def. 4. Choosing  $B \in \sigma(\mathbf{X})$  and  $u \in \mathcal{S}_U$ , we need to show that

$$\mathbf{E}[\mathbf{1}_B \nu_{\hat{\mathbf{X}}}(u)] = \mathbf{E}[\mathbf{1}_B \mathbf{1}_{\{u\}}(\tilde{\mathbf{U}})]. \quad (77)$$

The statement for  $B \in \sigma(\hat{\mathbf{X}})$  then follows by  $\sigma(\hat{\mathbf{X}}) \subseteq \sigma(\mathbf{X})$ , i.e., the  $\sigma(\mathbf{X})$ -measurability of  $\hat{\mathbf{X}}$ . We prove (77) by

$$\mathbf{E}[\mathbf{1}_B \nu_{\hat{\mathbf{X}}}(u)] = \sum_{i \in \mathcal{I}} \mathbf{E}[\mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_B \nu_i(u)] \quad (78)$$

$$= \sum_{i \in \mathcal{I}} \nu_i(u) \mathbf{E}[\mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_B] \quad (79)$$

$$= \sum_{i \in \mathcal{I}} \mathbf{E}[\mathbf{1}_u(\tilde{\mathbf{U}}_i)] \mathbf{E}[\mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_B] \quad (80)$$

$$= \sum_{i \in \mathcal{I}} \mathbf{E}[\mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}})] \quad (81)$$

$$= \sum_{i \in \mathcal{I}} \mathbf{E}[\mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}})] \quad (82)$$

$$= \mathbf{E}[\mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}})], \quad (83)$$

where we used Fubini's theorem [18, Thm. 14.16] in (81).

To prove  $\mathbf{I}(\mathbf{Y}; \tilde{\mathbf{U}}|\mathbf{X}) = 0$ , we need to show that for every  $y \in \mathcal{S}_Y$ ,  $u \in \mathcal{S}_U$ , and  $B \in \sigma(\mathbf{X})$ , we have

$$\int \mathbf{1}_B \kappa_{\mathbf{Y}|\mathbf{X}}(\cdot; y) \nu_{\hat{\mathbf{X}}}(u) d\mu_{\mathbf{YXU}} = \int \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}}) \mathbf{1}_y(\mathbf{Y}) d\mu_{\mathbf{YXU}\tilde{\mathbf{U}}_T} \quad (84)$$

and by integrating, we indeed obtain

$$\int \mathbf{1}_B \kappa_{\mathbf{Y}|\mathbf{X}}(\cdot; y) \nu_{\hat{\mathbf{X}}}(u) d\mu_{\mathbf{YXU}} \quad (85)$$

$$= \sum_{i \in \mathcal{I}} \int \mathbf{1}_B \mathbf{1}_i(\hat{\mathbf{X}}) \kappa_{\mathbf{Y}|\mathbf{X}}(\cdot; y) \nu_i(u) d\mu_{\mathbf{YXU}} \quad (86)$$

$$= \sum_{i \in \mathcal{I}} \nu_i(u) \int \mathbf{1}_B \mathbf{1}_i(\hat{\mathbf{X}}) \kappa_{\mathbf{Y}|\mathbf{X}}(\cdot; y) d\mu_{\mathbf{YXU}} \quad (87)$$

$$= \sum_{i \in \mathcal{I}} \int \mathbf{1}_u(\tilde{\mathbf{U}}_i) d\mu_{\tilde{\mathbf{U}}_T} \int \mathbf{1}_B \mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_y(\mathbf{Y}) d\mu_{\mathbf{YXU}} \quad (88)$$

$$= \sum_{i \in \mathcal{I}} \int \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}}_i) \mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_y(\mathbf{Y}) d\mu_{\mathbf{YXU}\tilde{\mathbf{U}}_T} \quad (89)$$

$$= \sum_{i \in \mathcal{I}} \int \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}}) \mathbf{1}_i(\hat{\mathbf{X}}) \mathbf{1}_y(\mathbf{Y}) d\mu_{\mathbf{YXU}\tilde{\mathbf{U}}_T} \quad (90)$$

$$= \int \mathbf{1}_B \mathbf{1}_u(\tilde{\mathbf{U}}) \mathbf{1}_y(\mathbf{Y}) d\mu_{\mathbf{YXU}\tilde{\mathbf{U}}_T}, \quad (91)$$

where we used part 2 of Def. 4 in (88) and Fubini's theorem [18, Thm. 14.16] in (89). By replacing  $\kappa_{\mathbf{Y}|\mathbf{X}}$  with  $\kappa_{\mathbf{Y}|\hat{\mathbf{X}}}$  and using  $B \in \sigma(\hat{\mathbf{X}})$ , the same argument can be used to show  $\mathbf{I}(\mathbf{Y}; \tilde{\mathbf{U}}|\hat{\mathbf{X}}) = 0$ .

## ACKNOWLEDGMENT

The authors would like to thank Michael Meidlinger for providing inspiration for this work.

## REFERENCES

- [1] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," in *Proc. 37<sup>th</sup> Annual Allerton Conference on Communication, Control, and Computing*, Monticello, IL, Sep. 1999, pp. 368–377.
- [2] N. Slonim and N. Tishby, "Document clustering using word clusters via the information bottleneck method," in *Proc. 23<sup>rd</sup> Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Athens, Greece, Jul. 2000, pp. 208–215.
- [3] S. Gordon, H. Greenspan, and J. Goldberger, "Applying the information bottleneck principle to unsupervised clustering of discrete and continuous image representations," in *Proc. Ninth IEEE International Conference on Computer Vision*, Nice, France, Oct. 2003, pp. 370–377.
- [4] G. Zeitler, R. Kötter, G. Bauch, and J. Widmer, "On quantizer design for soft values in the multiple-access relay channel," in *Proc. IEEE ICC*, Dresden, Germany, Jun. 2009.
- [5] G. Zeitler, A. Singer, and G. Kramer, "Low-precision A/D conversion for maximum information rate in channels with memory," *IEEE Trans. Communications*, vol. 60, no. 9, pp. 2511–2521, Sep. 2012.
- [6] A. Winkelbauer and G. Matz, "Joint network-channel coding for the asymmetric multiple-access relay channel," in *Proc. IEEE ICC*, Ottawa, Canada, Jun. 2012, pp. 2485–2489.
- [7] A. Winkelbauer, S. Farthofer, and G. Matz, "The rate-information trade-off for Gaussian vector channels," in *IEEE Int. Symp. Information Theory*, Honolulu, HI, USA, Jun. 2014, pp. 2849–2853.
- [8] A. Winkelbauer and G. Matz, "On quantization of log-likelihood ratios for maximum mutual information," in *Proc. IEEE SPAWC*, Stockholm, Sweden, Jun. 2015, pp. 316–320.
- [9] G. Pichler, P. Piantanida, and G. Matz, "Distributed information-theoretic biclustering of two memoryless sources," in *Proc. 53<sup>rd</sup> Annual Allerton Conference on Communication, Control, and Computing*, Monticello, IL, Sep. 2015, pp. 426–433.
- [10] 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.
- [11] G. Chechik, A. Globerson, N. Tishby, and Y. Weiss, "Information bottleneck for Gaussian variables," *Journal of Machine Learning Research*, vol. 6, pp. 165–188, Jan. 2005.
- [12] N. Slonim and N. Tishby, "Agglomerative information bottleneck," in *Proc. 12<sup>th</sup> Annu. Conf. on Neural Inform. Process. Syst.*, Denver, CO, 2000, pp. 617–623.
- [13] B. M. Kurkoski, "On the relationship between the KL means algorithm and the information bottleneck method," in *Proc. 11<sup>th</sup> International ITG Conference on Systems, Communications and Coding (SCC)*, Hamburg, Germany, Feb. 2017.
- [14] M. B. Westover and J. A. O'Sullivan, "Achievable rates for pattern recognition," *IEEE Trans. Inf. Theory*, vol. 54, no. 1, pp. 299–320, Jan. 2008.
- [15] R. Ahlswede and I. Csiszár, "Hypothesis testing with communication constraints," *IEEE Trans. Inf. Theory*, vol. 32, no. 4, pp. 533–542, Jul. 1986.
- [16] T. S. Han, "Hypothesis testing with multiterminal data compression," *IEEE Trans. Inf. Theory*, vol. 33, no. 6, pp. 759–772, Nov. 1987.
- [17] R. M. Gray, *Entropy and Information Theory*, 1st ed. Springer, 2013.
- [18] A. Klenke, *Probability Theory*. Springer, Sep. 2013.
- [19] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*. Cambridge University Press, 2011.