

# Vector Gaussian CEO Problem Under Logarithmic Loss and Applications

Yiğit Uğur <sup>†‡</sup>Iñaki Estella Aguerri <sup>†</sup>Abdellatif Zaidi <sup>†‡</sup>

<sup>†</sup> Mathematical and Algorithmic Sciences Lab, Paris Research Center, Huawei Technologies,  
Boulogne-Billancourt, 92100, France

<sup>‡</sup> Université Paris-Est, Champs-sur-Marne, 77454, France

{yigit.ugur@gmail.com, inaki.estella@gmail.com, abdellatif.zaidi@u-pem.fr}

## Abstract

In this paper, we study the vector Gaussian Chief Executive Officer (CEO) problem under logarithmic loss distortion measure. Specifically,  $K \geq 2$  agents observe independently corrupted Gaussian noisy versions of a remote vector Gaussian source, and communicate independently with a decoder or CEO over rate-constrained noise-free links. The CEO also has its own Gaussian noisy observation of the source and wants to reconstruct the remote source to within some prescribed distortion level where the incurred distortion is measured under the logarithmic loss penalty criterion. We find an explicit characterization of the rate-distortion region of this model. The result can be seen as the counterpart to the vector Gaussian setting of that by Courtade-Weissman which provides the rate-distortion region of the model in the discrete memoryless setting. For the proof of this result, we obtain an outer bound by means of a technique that relies on the de Bruijn identity and the properties of Fisher information. The approach is similar to Ekrem-Ululuk outer bounding technique for the vector Gaussian CEO problem under quadratic distortion measure, for which it was there found generally non-tight; but it is shown here to yield a complete characterization of the region for the case of logarithmic loss measure. Also, we show that Gaussian test channels with time-sharing exhaust the Berger-Tung inner bound, which is optimal. Furthermore, application of our results allows us to find the complete solutions of two related problems: a quadratic vector Gaussian CEO problem with *determinant* constraint and the vector Gaussian distributed Information Bottleneck problem. Finally, we develop Blahut-Arimoto type algorithms that allow to compute numerically the regions provided in this paper, for both discrete and Gaussian models. With the known relevance of the logarithmic loss fidelity measure in the context of learning and prediction, the proposed algorithms may find usefulness in a variety of applications where learning is performed distributively. We illustrate the efficiency of our algorithms through some numerical examples.

The results of this paper have been presented in part at the 2017 IEEE Information Theory Workshop [1] and in part at the 2018 IEEE Information Theory Workshop [2].

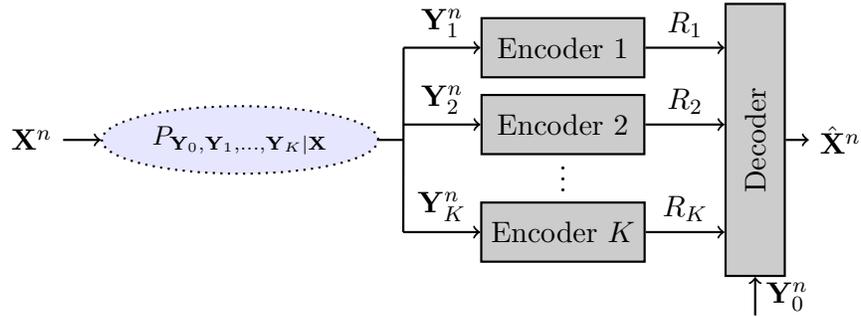


Fig. 1: Chief Executive Officer (CEO) source coding problem with side information.

## I. INTRODUCTION

Consider the vector Gaussian Chief Executive Officer (CEO) problem shown in Figure 1. In this model, there are  $K \geq 2$  agents each observing a noisy version of a vector Gaussian source  $\mathbf{X}$ . The goal of the agents is to describe the source to a central unit, which wants to reconstruct this source to within a prescribed distortion level. The incurred distortion is measured according to some loss measure  $d : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathbb{R}$ , where  $\hat{\mathcal{X}}$  designates the reconstruction alphabet. For quadratic distortion measure, i.e.,

$$d(x, \hat{x}) = |x - \hat{x}|^2,$$

the rate-distortion region of the vector Gaussian CEO problem is still unknown in general, except in few special cases the most important of which is perhaps the case of scalar sources, i.e., scalar Gaussian CEO problem. For this case, a complete solution, in terms of characterization of the optimal rate-distortion region, was found independently by Oohama in [3] and by Prabhakaran *et al.* in [4]. Key to establishing this result is a judicious application of the entropy power inequality. The extension of this argument to the case of vector Gaussian sources, however, is not straightforward as the entropy power inequality is known to be non-tight in this setting. The reader may refer also to [5], [6] where non-tight outer bounds on the rate-distortion region of the vector Gaussian CEO problem under quadratic distortion measure are obtained by establishing some extremal inequalities that are similar to Liu-Viswanath [7], and to [8] where a strengthened extremal inequality yields a complete characterization of the region of the vector Gaussian CEO problem in the special case of trace distortion constraint.

In this paper, we study the CEO problem of Figure 1 in the case in which  $(\mathbf{X}, \mathbf{Y}_0, \mathbf{Y}_1, \dots, \mathbf{Y}_K)$  is jointly Gaussian and the distortion is measured using the logarithmic loss criterion, i.e.,

$$d^{(n)}(x^n, \hat{x}^n) = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i), \quad (1)$$

with the letter-wise distortion given by

$$d(x, \hat{x}) = \log \frac{1}{\hat{x}(x)}, \quad (2)$$

where  $\hat{x}(\cdot)$  designates a probability distribution on  $\mathcal{X}$  and  $\hat{x}(x)$  is the value of this distribution evaluated for the outcome  $x \in \mathcal{X}$ .

The logarithmic loss distortion measure, often referred to as *self-information loss* in the literature about prediction, plays a central role in settings in which reconstructions are allowed to be ‘soft’, rather than ‘hard’ or deterministic. That is, rather than just assigning a deterministic value to each sample of the source, the decoder also gives an assessment of the degree of confidence or reliability on each estimate, in the form of weights or probabilities. This measure, which was introduced in the context of rate-distortion theory by Courtade *et al.* [9], [10], has appreciable mathematical properties [11], [12], such as a deep connection to lossless coding for which fundamental limits are well developed (e.g., see [13] for recent results on universal lossy compression under logarithmic loss that are built on this connection). Also, it is widely used as a penalty criterion in various contexts, including clustering and classification [14], pattern recognition, learning and prediction [15], image processing [16], secrecy [17] and others.

#### A. Main Contributions

The main contribution of this paper is a complete characterization of the rate-distortion region of the vector Gaussian CEO problem of Figure 1 under logarithmic loss distortion measure. The result can be seen as the counterpart to the vector Gaussian case of that by Courtade and Weissman [10, Theorem 10], who established the rate-distortion region of the CEO problem under logarithmic loss in the discrete memoryless (DM) case. For the proof of this result, we derive a matching outer bound by means of a technique that relies of the de Bruijn identity, a connection between differential entropy and Fisher information, along with the properties of minimum mean square error (MMSE) and Fisher information. By opposition to the case of quadratic distortion measure, for which the application of this technique was shown in [18] to result in an outer bound that is generally non-tight, we show that this approach is successful in the case of logarithmic distortion measure and yields a complete characterization of the region. On this aspect, it is noteworthy that, in the specific case of scalar Gaussian sources, an alternate converse proof may be obtained by extending that of the scalar Gaussian many-help-one source coding problem by Oahama [3] and Prabhakaran *et al.* [4] by accounting for side information and replacing the original mean square error distortion constraint with conditional entropy. However, such approach does not seem to lead to a conclusive result in the vector case as the entropy power inequality is known to be generally non-tight in this setting [19], [20]. The proof of the achievability part simply follows by evaluating a straightforward extension to the continuous alphabet case of the solution of the DM model using Gaussian test channels and *no* time-sharing. Because this does *not* necessarily imply that Gaussian test channels also exhaust the Berger-Tung inner bound, we investigate the question and we show that they *do* if time-sharing is allowed.

Furthermore, we show that application of our results allows us to find complete solutions to two related problems. The first is a quadratic vector Gaussian CEO problem with reconstruction constraint on the *determinant* of the error covariance matrix that we introduce here, and for which we also characterize the optimal rate-distortion region. Key to establishing this result, we show that the rate-distortion region of vector Gaussian CEO problem under logarithmic loss which is found in this paper translates into an outer bound on the rate region of the quadratic

vector Gaussian CEO problem with *determinant* constraint. The reader may refer to, e.g., [21] and [22] for examples of usage of such a determinant constraint in the context of equalization and others. The second is an extension of Tishby’s single-encoder Information Bottleneck (IB) method [14] to the case of multiple encoders. Information theoretically, this problem is known to be essentially a remote source coding problem with logarithmic loss distortion measure [23]; and, so, we use our result for the vector Gaussian CEO problem under logarithmic loss to infer a full characterization of the optimal trade-off between *complexity* (or rate) and *accuracy* (or information) for the distributed vector Gaussian IB problem.

Finally, for both DM and memoryless Gaussian settings we develop Blahut-Arimoto (BA) [24], [25] type iterative algorithms that allow to compute (approximations of) the rate regions that are established in this paper; and prove their convergence to stationary points. We do so through a variational formulation that allows to determine the set of self-consistent equations that are satisfied by the stationary solutions. In the Gaussian case, we show that the algorithm reduces to an appropriate updating rule of the parameters of noisy linear projections. We note that the computation of the rate-distortion regions of multiterminal and CEO source coding problems is important *per-se* as it involves non-trivial optimization problems over distributions of auxiliary random variables. Also, since the logarithmic loss function is instrumental in connecting problems of multiterminal rate-distortion theory with those of distributed learning and estimation, the algorithms that are developed in this paper also find usefulness in emerging applications in those areas. For example, our algorithm for the DM CEO problem under logarithm loss measure can be seen as a generalization of Tishby’s IB method [14] to the distributed learning setting. Similarly, our algorithm for the vector Gaussian CEO problem under logarithm loss measure can be seen as a generalization of that of [26]–[28] to the distributed learning setting. For other extensions of the BA algorithm in the context of multiterminal data transmission and compression, the reader may refer to related works on point-to-point [29], [30] and broadcast and multiple access multiterminal settings [31], [32].

### B. Related Works

As we already mentioned, this paper mostly relates to [10] in which the authors establish the rate-distortion region of the DM CEO problem under logarithmic loss in the case of an arbitrary number of encoders and no side information at the decoder, as well as that of the DM multiterminal source coding problem under logarithmic loss in the case of two encoders and no side information at the decoder. Motivated by the increasing interest for problems of learning and prediction, a growing body of works study point-to-point and multiterminal source coding models under logarithmic loss. In [11], Jiao *et al.* provide a fundamental justification for inference using logarithmic loss, by showing that under some mild conditions (the loss function satisfying some data processing property and alphabet size larger than two) the reduction in optimal risk in the presence of side information is uniquely characterized by mutual information, and the corresponding loss function coincides with the logarithmic loss. Somewhat related, in [33] Painsky and Wornell show that for binary classification problems the logarithmic loss dominates “universally” any other convenient (i.e., smooth, proper and convex) loss function, in the sense that by minimizing the logarithmic loss one minimizes the regret that is associated with any such measures. More

specifically, the divergence associated any smooth, proper and convex loss function is shown to be bounded from above by the Kullback-Leibler divergence, up to a multiplicative normalization constant. In [13], the authors study the problem of universal lossy compression under logarithmic loss, and derive bounds on the non-asymptotic fundamental limit of fixed-length universal coding with respect to a family of distributions that generalize the well-known minimax bounds for universal lossless source coding. In [34], the minimax approach is studied for a problem of remote prediction and is shown to correspond to a one-shot minimax noisy source coding problem. The setting of remote prediction of [34] provides an approximate one-shot operational interpretation of the Information Bottleneck method of [14], which is also sometimes interpreted as a remote source coding problem under logarithmic loss [23].

Logarithmic loss is also instrumental in problems of data compression under a mutual information constraint [35], and problems of relaying with relay nodes that are constrained not to know the users' codebooks (sometimes termed "oblivious" or nomadic processing) which is studied in the single user case first by Sanderovich *et al.* in [36] and then by Simeone *et al.* in [37], and in the multiple user multiple relay case by Aguerri *et al.* in [38] and [39]. Other applications in which the logarithmic loss function can be used include secrecy and privacy [17], [40], hypothesis testing against independence [41]–[45] and others.

### C. Outline and Notation

The rest of this paper is organized as follows. Section II provides a formal description of the vector Gaussian CEO model that we study in this paper, as well as some definitions that are related to it. Section III contains the main results of this paper: an explicit characterization of the rate-distortion region of the memoryless vector Gaussian CEO problem with side information under logarithmic loss as well as the proof that Gaussian test channels with time-sharing exhaust the Berger-Tung rate region which is optimal. In Section IV we use our results on the CEO problem under logarithmic loss to infer complete solutions of two related problems: a quadratic vector Gaussian CEO problem with a determinant constraint on the covariance matrix error and the vector Gaussian distributed Information Bottleneck problem. Section V provides BA-type algorithms for the computation of the rate-distortion regions that are established in this paper in both DM and Gaussian cases as well as proofs of their convergence and some numerical examples. The proofs are deferred to the appendices section.

Throughout this paper, we use the following notation. Upper case letters are used to denote random variables, e.g.,  $X$ ; lower case letters are used to denote realizations of random variables, e.g.,  $x$ ; and calligraphic letters denote sets, e.g.,  $\mathcal{X}$ . The cardinality of a set  $\mathcal{X}$  is denoted by  $|\mathcal{X}|$ . The length- $n$  sequence  $(X_1, \dots, X_n)$  is denoted as  $X^n$ ; and, for integers  $j$  and  $k$  such that  $1 \leq k \leq j \leq n$ , the sub-sequence  $(X_k, X_{k+1}, \dots, X_j)$  is denoted as  $X_k^j$ . Probability mass functions (pmfs) are denoted by  $P_X(x) = \Pr\{X = x\}$ ; and, sometimes, for short, as  $p(x)$ . We use  $\mathcal{P}(\mathcal{X})$  to denote the set of discrete probability distributions on  $\mathcal{X}$ . Boldface upper case letters denote vectors or matrices, e.g.,  $\mathbf{X}$ , where context should make the distinction clear. For an integer  $K \geq 1$ , we denote the set of integers smaller or equal  $K$  as  $\mathcal{K} = \{k \in \mathbb{N} : 1 \leq k \leq K\}$ . For a set of integers  $\mathcal{S} \subseteq \mathcal{K}$ , the complementary set of  $\mathcal{S}$  is denoted by  $\mathcal{S}^c$ , i.e.,  $\mathcal{S}^c = \{k \in \mathbb{N} : k \in \mathcal{K} \setminus \mathcal{S}\}$ . Sometimes, for convenience we will need to

define  $\bar{\mathcal{S}}$  as  $\bar{\mathcal{S}} = \{0\} \cup \mathcal{S}^c$ . For a set of integers  $\mathcal{S} \subseteq \mathcal{K}$ ; the notation  $X_{\mathcal{S}}$  designates the set of random variables  $\{X_k\}$  with indices in the set  $\mathcal{S}$ , i.e.,  $X_{\mathcal{S}} = \{X_k\}_{k \in \mathcal{S}}$ . We denote the covariance of a zero mean, complex-valued, vector  $\mathbf{X}$  by  $\Sigma_{\mathbf{x}} = \mathbb{E}[\mathbf{X}\mathbf{X}^\dagger]$ , where  $(\cdot)^\dagger$  indicates conjugate transpose. Similarly, we denote the cross-correlation of two zero-mean vectors  $\mathbf{X}$  and  $\mathbf{Y}$  as  $\Sigma_{\mathbf{x},\mathbf{y}} = \mathbb{E}[\mathbf{X}\mathbf{Y}^\dagger]$ , and the conditional correlation matrix of  $\mathbf{X}$  given  $\mathbf{Y}$  as  $\Sigma_{\mathbf{x}|\mathbf{y}} = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}])(\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}])^\dagger]$  i.e.,  $\Sigma_{\mathbf{x}|\mathbf{y}} = \Sigma_{\mathbf{x}} - \Sigma_{\mathbf{x},\mathbf{y}}\Sigma_{\mathbf{y}}^{-1}\Sigma_{\mathbf{y},\mathbf{x}}$ . For matrices  $\mathbf{A}$  and  $\mathbf{B}$ , the notation  $\text{diag}(\mathbf{A}, \mathbf{B})$  denotes the block diagonal matrix whose diagonal elements are the matrices  $\mathbf{A}$  and  $\mathbf{B}$  and its off-diagonal elements are the all zero matrices. Also, for a set of integers  $\mathcal{J} \subset \mathbb{N}$  and a family of matrices  $\{\mathbf{A}_i\}_{i \in \mathcal{J}}$  of the same size, the notation  $\mathbf{A}_{\mathcal{J}}$  is used to denote the (super) matrix obtained by concatenating vertically the matrices  $\{\mathbf{A}_i\}_{i \in \mathcal{J}}$ , where the indices are sorted in the ascending order, e.g,  $\mathbf{A}_{\{0,2\}} = [\mathbf{A}_0^\dagger, \mathbf{A}_2^\dagger]^\dagger$ .

## II. PROBLEM FORMULATION

Consider the  $K$ -encoder CEO problem with side information shown in Figure 1, where the agents observations are assumed to be Gaussian noisy versions of a remote vector Gaussian source. Specifically, let  $(\mathbf{X}, \mathbf{Y}_0, \mathbf{Y}_1, \dots, \mathbf{Y}_K)$  be a jointly Gaussian random vector, with zero mean and covariance matrix  $\Sigma_{(\mathbf{x}, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_K)}$ . Without loss of generality, the remote vector source  $\mathbf{X} \in \mathbb{C}^{n_x}$  is assumed complex-valued, has  $n_x$ -dimensions, and is assumed to be Gaussian with zero mean and covariance matrix  $\Sigma_{\mathbf{x}} \succeq \mathbf{0}$ .  $\mathbf{X}^n = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  denotes a collection of  $n$  independent copies of  $\mathbf{X}$ . The agents' observations are Gaussian noisy versions of the remote vector source, with the observation at agent  $k \in \mathcal{K}$  given by

$$\mathbf{Y}_{k,i} = \mathbf{H}_k \mathbf{X}_i + \mathbf{N}_{k,i}, \quad \text{for } i = 1, \dots, n, \quad (3)$$

where  $\mathbf{H}_k \in \mathbb{C}^{n_k \times n_x}$  represents the channel matrix connecting the remote vector source to the  $k$ -th agent; and  $\mathbf{N}_{k,i} \in \mathbb{C}^{n_k}$  is the noise vector at this agent, assumed to be i.i.d. Gaussian with zero-mean and independent from  $\mathbf{X}_i$ . The decoder has its own noisy observation of the remote vector source, in the form of a correlated jointly Gaussian side information stream  $\mathbf{Y}_0^n$ , with

$$\mathbf{Y}_{0,i} = \mathbf{H}_0 \mathbf{X}_i + \mathbf{N}_{0,i}, \quad \text{for } i = 1, \dots, n, \quad (4)$$

where, similar to the above,  $\mathbf{H}_0 \in \mathbb{C}^{n_0 \times n_x}$  is the channel matrix connecting the remote vector source to the CEO; and  $\mathbf{N}_{0,i} \in \mathbb{C}^{n_0}$  is the noise vector at the CEO, assumed to be Gaussian with zero-mean and covariance matrix  $\Sigma_0 \succeq \mathbf{0}$  and independent from  $\mathbf{X}_i$ . In this section, it is assumed that the agents' observations are independent conditionally given the remote vector source  $\mathbf{X}^n$  and the side information  $\mathbf{Y}_0^n$ , i.e., for all  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$\mathbf{Y}_{\mathcal{S}}^n \text{---} (\mathbf{X}^n, \mathbf{Y}_0^n) \text{---} \mathbf{Y}_{\mathcal{S}^c}^n. \quad (5)$$

Using (3) and (4), it is easy to see that the assumption (5) is equivalent to that the noises at the agents are independent conditionally given  $\mathbf{N}_0$ . For notational simplicity,  $\Sigma_k$  denotes the conditional covariance matrix of the noise  $\mathbf{N}_k$  at the  $k$ -th agent given  $\mathbf{N}_0$ , i.e.,  $\Sigma_k := \Sigma_{\mathbf{n}_k | \mathbf{n}_0}$ . Recalling that for a set  $\mathcal{S} \subseteq \mathcal{K}$ ,  $\mathbf{N}_{\mathcal{S}}$  designates the collection of noise vectors with indices in the set  $\mathcal{S}$ , in what follows we denote the covariance matrix of  $\mathbf{N}_{\mathcal{S}}$  as  $\Sigma_{\mathbf{n}_{\mathcal{S}}}$ .

In this model, Encoder (or agent)  $k \in \mathcal{K}$  uses  $R_k$  bits per sample to describe its observation  $\mathbf{Y}_k^n$  to the decoder. The decoder wants to reconstruct the remote source  $\mathbf{X}^n$  to within a prescribed fidelity level. Similar to [10], we consider the reproduction alphabet to be equal to the set of probability distributions over the source alphabet  $\mathbb{C}^{n \times n_x}$ . In other words, the decoder generates ‘soft’ estimates of the remote source’s sequences. We consider the logarithmic loss distortion measure defined as in (1), where the letter-wise distortion measure is given by (2).

**Definition 1.** A rate-distortion code (of blocklength  $n$ ) for the model of Figure 1 consists of  $K$  encoding functions

$$\phi_k^{(n)} : \mathbb{C}^{n \times n_k} \rightarrow \{1, \dots, M_k^{(n)}\}, \quad \text{for } k = 1, \dots, K,$$

and a decoding function

$$\psi^{(n)} : \{1, \dots, M_1^{(n)}\} \times \dots \times \{1, \dots, M_K^{(n)}\} \times \mathbb{C}^{n \times n_0} \rightarrow \mathcal{P}(\mathbb{C}^{n \times n_x}),$$

where  $\mathcal{P}(\mathbb{C}^{n \times n_x})$  designates the set of probability distributions over the  $n$ -Cartesian product of  $\mathbb{C}^{n_x}$ . ■

**Definition 2.** A rate-distortion tuple  $(R_1, \dots, R_K, D)$  is achievable for the vector Gaussian CEO problem with side information if there exist a blocklength  $n$ ,  $K$  encoding functions  $\{\phi_k^{(n)}\}_{k=1}^K$  and a decoding function  $\psi^{(n)}$  such that

$$\begin{aligned} R_k &\geq \frac{1}{n} \log M_k^{(n)}, \quad \text{for } k = 1, \dots, K, \\ D &\geq \mathbb{E}[d^{(n)}(\mathbf{X}^n, \psi^{(n)}(\phi_1^{(n)}(\mathbf{Y}_1^n), \dots, \phi_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n))]. \end{aligned}$$

The rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^*$  of the vector Gaussian CEO problem under logarithmic loss is defined as the union of all non-negative tuples  $(R_1, \dots, R_K, D)$  that are achievable. ■

The main goal of this paper is to characterize the rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^*$  of the vector Gaussian CEO problem under logarithmic loss.

### III. MAIN RESULTS

In this section we provide an explicit characterization of the rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^*$  of the vector Gaussian CEO problem under logarithmic loss. Also, we show that Gaussian test channels with time-sharing exhaust the Berger-Tung region which is optimal.

#### A. Rate-Distortion Region

We first state the following theorem which follows essentially by an easy application of [10, Theorem 10] that provides the rate-distortion region of the DM version of the problem.

**Definition 3.** For given tuple of auxiliary random variables  $(U_1, \dots, U_K, Q)$  with distribution  $P_{U_{\mathcal{K}}, Q}(u_{\mathcal{K}}, q)$  such that  $P_{\mathbf{X}, \mathbf{Y}_0, \mathbf{Y}_{\mathcal{K}}, U_{\mathcal{K}}, Q}(\mathbf{x}, \mathbf{y}_0, \mathbf{y}_{\mathcal{K}}, u_{\mathcal{K}}, q)$  factorizes as

$$P_{\mathbf{X}, \mathbf{Y}_0}(\mathbf{x}, \mathbf{y}_0) \prod_{k=1}^K P_{\mathbf{Y}_k | \mathbf{X}, \mathbf{Y}_0}(\mathbf{y}_k | \mathbf{x}, \mathbf{y}_0) P_Q(q) \prod_{k=1}^K P_{U_k | \mathbf{Y}_k, Q}(u_k | \mathbf{y}_k, q), \quad (6)$$

define  $\mathcal{RD}_{\text{CEO}}^{\text{I}}(U_1, \dots, U_K, Q)$  as the set of all non-negative rate-distortion tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$D + \sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} I(\mathbf{Y}_k; U_k | \mathbf{X}, \mathbf{Y}_0, Q) + h(\mathbf{X} | U_{\mathcal{S}^c}, \mathbf{Y}_0, Q). \quad (7)$$

Also, let  $\mathcal{RD}_{\text{CEO}}^{\text{I}} := \bigcup \mathcal{RD}_{\text{CEO}}^{\text{I}}(U_1, \dots, U_K, Q)$  where the union is taken over all tuples  $(U_1, \dots, U_K, Q)$  with distributions that satisfy (6). ■

**Definition 4.** For given tuple of auxiliary random variables  $(V_1, \dots, V_K, Q')$  with distribution  $P_{V_{\mathcal{K}}, Q'}(v_{\mathcal{K}}, q')$  such that  $P_{\mathbf{X}, \mathbf{Y}_0, \mathbf{Y}_{\mathcal{K}}, V_{\mathcal{K}}, Q'}(\mathbf{x}, \mathbf{y}_0, \mathbf{y}_{\mathcal{K}}, v_{\mathcal{K}}, q')$  factorizes as

$$P_{\mathbf{X}, \mathbf{Y}_0}(\mathbf{x}, \mathbf{y}_0) \prod_{k=1}^K P_{\mathbf{Y}_k | \mathbf{X}, \mathbf{Y}_0}(\mathbf{y}_k | \mathbf{x}, \mathbf{y}_0) P_{Q'}(q') \prod_{k=1}^K P_{V_k | \mathbf{Y}_k, Q'}(v_k | \mathbf{y}_k, q'), \quad (8)$$

define  $\mathcal{RD}_{\text{CEO}}^{\text{II}}(V_1, \dots, V_K, Q')$  as the set of all non-negative rate-distortion tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$\sum_{k \in \mathcal{S}} R_k \geq I(\mathbf{Y}_{\mathcal{S}}; V_{\mathcal{S}} | V_{\mathcal{S}^c}, \mathbf{Y}_0, Q') \quad (9)$$

$$D \geq h(\mathbf{X} | V_1, \dots, V_K, \mathbf{Y}_0, Q'). \quad (10)$$

Also, let  $\mathcal{RD}_{\text{CEO}}^{\text{II}} := \bigcup \mathcal{RD}_{\text{CEO}}^{\text{II}}(V_1, \dots, V_K, Q')$  where the union is taken over all tuples  $(V_1, \dots, V_K, Q')$  with distributions that satisfy (8). ■

**Theorem 1.** The rate-distortion region for the vector Gaussian CEO problem under logarithmic loss is given by

$$\mathcal{RD}_{\text{VG-CEO}}^* = \mathcal{RD}_{\text{CEO}}^{\text{I}} = \mathcal{RD}_{\text{CEO}}^{\text{II}}.$$

*Proof.* The proof of Theorem 1 is given in Appendix I. □

For convenience, we now introduce the following notation which will be instrumental in what follows. Let, for every set  $\mathcal{S} \subseteq \mathcal{K}$ , the set  $\bar{\mathcal{S}} = \{0\} \cup \mathcal{S}^c$ . Also, for  $\mathcal{S} \subseteq \mathcal{K}$  and given matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ , let  $\mathbf{\Lambda}_{\bar{\mathcal{S}}}$  designate the block-diagonal matrix given by

$$\mathbf{\Lambda}_{\bar{\mathcal{S}}} := \begin{bmatrix} \mathbf{0} & & \mathbf{0} \\ \mathbf{0} & \text{diag}(\{\mathbf{\Sigma}_k - \mathbf{\Sigma}_k \mathbf{\Omega}_k \mathbf{\Sigma}_k\}_{k \in \mathcal{S}^c}) & \end{bmatrix}, \quad (11)$$

where  $\mathbf{0}$  in the principal diagonal elements is the  $n_0 \times n_0$ -all zero matrix. Besides, the notation  $\mathbf{H}_{\mathcal{S}}$  is used to denote the (super) matrix obtained by concatenating vertically the matrices  $\{\mathbf{H}_i\}_{i \in \mathcal{S}}$ , where the indices are sorted in the ascending order, e.g,  $\mathbf{H}_{\{0,2\}} = [\mathbf{H}_0^\dagger, \mathbf{H}_2^\dagger]^\dagger$ .

The following theorem is the main contribution of this paper, which is an explicit characterization of the rate-distortion region of the vector Gaussian CEO problem with side information under logarithmic loss measure.

**Theorem 2.** *The rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^*$  of the vector Gaussian CEO problem under logarithmic loss is given by the set of all non-negative rate-distortion tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,*

$$D + \sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} \log \frac{1}{|\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k|} + \log \left| (\pi e) \left( \mathbf{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{H}_{\bar{\mathcal{S}}}^\dagger \mathbf{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1} (\mathbf{I} - \mathbf{\Lambda}_{\bar{\mathcal{S}}} \mathbf{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1}) \mathbf{H}_{\bar{\mathcal{S}}} \right)^{-1} \right|,$$

for matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ , where  $\bar{\mathcal{S}} = \{0\} \cup \mathcal{S}^c$  and  $\mathbf{\Lambda}_{\bar{\mathcal{S}}}$  is as defined by (11).

*Proof.* The proof of the direct part of Theorem 2 follows simply by evaluating the region  $\mathcal{RD}_{\text{CEO}}^{\text{I}}$  as described in Theorem 1 using Gaussian test channels and no time-sharing. Specifically, we set  $Q = \emptyset$  and  $p(u_k | \mathbf{y}_k, q) = \mathcal{CN}(\mathbf{y}_k, \mathbf{\Sigma}_k^{1/2} (\mathbf{\Omega}_k - \mathbf{I}) \mathbf{\Sigma}_k^{1/2})$ ,  $k \in \mathcal{K}$ . The proof of the converse appears in Appendix II.  $\square$

In the case in which the noises at the agents are independent among them and from the noise  $\mathbf{N}_0$  at the CEO, the result of Theorem 2 takes a simpler form which is stated in the following corollary.

**Corollary 1.** *Consider the vector Gaussian CEO problem described by (3) and (4) with the noises  $(\mathbf{N}_1, \dots, \mathbf{N}_K)$  being independent among them and with  $\mathbf{N}_0$ . Under logarithmic loss, the rate-distortion region this model is given by the set of all non-negative tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,*

$$D + \sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} \log \frac{1}{|\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k|} + \log \left| (\pi e) \left( \mathbf{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{H}_0^\dagger \mathbf{\Sigma}_0^{-1} \mathbf{H}_0 + \sum_{k \in \mathcal{S}^c} \mathbf{H}_k^\dagger \mathbf{\Omega}_k \mathbf{H}_k \right)^{-1} \right|,$$

for some matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ .  $\blacksquare$

**Remark 1.** *The direct part of Theorem 2 shows that Gaussian test channels and no-time sharing exhaust the region. For the converse proof of Theorem 2, we derive an outer bound on the region  $\mathcal{RD}_{\text{CEO}}^{\text{I}}$ . In doing so, we use the de Bruijn identity, a connection between differential entropy and Fisher information, along with the properties of MMSE and Fisher information. By opposition to the case of quadratic distortion measure for which the application of this technique was shown in [18] to result in an outer bound that is generally non-tight, Theorem 2 shows that the approach is successful in the case of logarithmic loss distortion measure as it yields a complete characterization of the region. On this aspect, note that in the specific case of scalar Gaussian sources, an alternate converse proof may be obtained by extending that of the scalar Gaussian many-help-one source coding problem by Oahama [3] and Prabhakaran et al. [4] through accounting for additional side information at CEO and replacing the original mean square error distortion constraint with conditional entropy. However, such approach does not seem conclusive in the vector case, as the entropy power inequality is known to be generally non-tight in this setting [19], [20].  $\blacksquare$*

**Remark 2.** *The result of Theorem 2 generalizes that of [35], which considers the case of only one agent, i.e., the remote vector Gaussian Wyner-Ziv model under logarithmic loss, to the case of an arbitrarily number of agents. The converse proof of [35], which relies on the technique of orthogonal transform to reduce the vector setting to one of parallel scalar Gaussian settings, seems insufficient to diagonalize all the noise covariance matrices simultaneously in the case of more than one agent. The result of Theorem 2 is also connected to recent developments on characterizing the capacity of multiple-input multiple-output (MIMO) relay channels in which the relay nodes are connected to the receiver through error-free finite-capacity links (i.e., the so-called cloud radio access networks). In particular, the reader may refer to [46, Theorem 4] where important progress is done, and [39] where compress-and-forward with joint decompression-decoding is shown to be optimal under the constraint of oblivious relay processing. ■*

### B. Gaussian Test Channels with Time-Sharing Exhaust the Berger-Tung Region

Theorem 1 shows that the union of all rate-distortion tuples that satisfy (7) for all subsets  $S \subseteq \mathcal{K}$  coincides with the Berger-Tung inner bound in which time-sharing is used. The direct part of Theorem 2 is obtained by evaluating (7) using Gaussian test channels and no time-sharing, i.e.,  $Q = \emptyset$ , not the Berger-Tung inner bound. The reader may wonder: i) whether Gaussian test channels also exhaust the Berger-Tung inner bound for the vector Gaussian CEO problem that we study here, and ii) whether time-sharing is needed with the Berger-Tung scheme. In this section, we answer both questions in the affirmative. In particular, we show that the Berger-Tung coding scheme with Gaussian test channels and time-sharing achieves distortion levels that are not larger than any other coding scheme.

**Proposition 1.** *The rate-distortion region for the vector Gaussian CEO problem under logarithmic loss is given by*

$$\mathcal{RD}_{\text{VG-CEO}}^* = \bigcup \mathcal{RD}_{\text{CEO}}^{\text{II}}(V_1^{\text{G}}, \dots, V_K^{\text{G}}, Q'),$$

where  $\mathcal{RD}_{\text{CEO}}^{\text{II}}(\cdot)$  is as given in Definition 4 and the superscript G is used to denote that the union is taken over Gaussian distributed  $V_k^{\text{G}} \sim p(v_k | \mathbf{y}_k, q')$  conditionally on  $(\mathbf{Y}_k, Q')$ .

*Proof.* For the proof of Proposition 1, it is sufficient to show that, for fixed Gaussian conditional distributions  $\{p(u_k | \mathbf{y}_k)\}_{k=1}^K$ , the extreme points of the polytopes defined by (7) are *dominated* by points that are in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$  and which are achievable using Gaussian conditional distributions  $\{p(v_k | \mathbf{y}_k, q')\}_{k=1}^K$ . Hereafter, we give a brief outline of proof for the case  $K = 2$ . The proof for  $K \geq 2$  follows similarly; and is omitted for brevity. Consider the inequalities (7) with  $Q = \emptyset$  and  $(U_1, U_2) := (U_1^{\text{G}}, U_2^{\text{G}})$  chosen to be Gaussian (see Theorem 2). Consider now the extreme points of the polytopes defined by the obtained inequalities:

$$\begin{aligned} P_1 &= (0, 0, I(\mathbf{Y}_1; U_1^{\text{G}} | \mathbf{X}, \mathbf{Y}_0) + I(\mathbf{Y}_2; U_2^{\text{G}} | \mathbf{X}, \mathbf{Y}_0) + h(\mathbf{X} | \mathbf{Y}_0)) \\ P_2 &= (I(\mathbf{Y}_1; U_1^{\text{G}} | \mathbf{Y}_0), 0, I(U_2^{\text{G}}; \mathbf{Y}_2 | \mathbf{X}, \mathbf{Y}_0) + h(\mathbf{X} | U_1^{\text{G}}, \mathbf{Y}_0)) \end{aligned}$$

$$\begin{aligned}
P_3 &= (0, I(\mathbf{Y}_2; U_2^G | \mathbf{Y}_0), I(U_1^G; \mathbf{Y}_1 | \mathbf{X}, \mathbf{Y}_0) + h(\mathbf{X} | U_2^G, \mathbf{Y}_0)) \\
P_4 &= (I(\mathbf{Y}_1; U_1^G | \mathbf{Y}_0), I(\mathbf{Y}_2; U_2^G | U_1^G, \mathbf{Y}_0), h(\mathbf{X} | U_1^G, U_2^G, \mathbf{Y}_0)) \\
P_5 &= (I(\mathbf{Y}_1; U_1^G | U_2^G, \mathbf{Y}_0), I(\mathbf{Y}_2; U_2^G | \mathbf{Y}_0), h(\mathbf{X} | U_1^G, U_2^G, \mathbf{Y}_0)),
\end{aligned}$$

where the point  $P_j$  is a triple  $(R_1^{(j)}, R_2^{(j)}, D^{(j)})$ . It is easy to see that each of these points is *dominated* by a point in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$ , i.e., there exists  $(R_1, R_2, D) \in \mathcal{RD}_{\text{CEO}}^{\text{II}}$  for which  $R_1 \leq R_1^{(j)}$ ,  $R_2 \leq R_2^{(j)}$  and  $D \leq D^{(j)}$ . To see this, first note that  $P_4$  and  $P_5$  are both in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$ . Next, observe that the point  $(0, 0, h(\mathbf{X} | \mathbf{Y}_0))$  is in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$ , which is clearly achievable by letting  $(V_1, V_2, Q') = (\emptyset, \emptyset, \emptyset)$ , dominates  $P_1$ . Also, by using letting  $(V_1, V_2, Q') = (U_1^G, \emptyset, \emptyset)$ , we have that the point  $(I(\mathbf{Y}_1; U_1^G | \mathbf{Y}_0), 0, h(\mathbf{X} | U_1^G, \mathbf{Y}_0))$  is in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$ , and dominates the point  $P_2$ . A similar argument shows that  $P_3$  is dominated by a point in  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$ . The proof is terminated by observing that, for all above corner points,  $V_k$  is set either equal  $U_k^G$  (which is Gaussian distributed conditionally on  $\mathbf{Y}_k$ ) or a constant.  $\square$

**Remark 3.** Proposition 1 shows that for the vector Gaussian CEO problem with side information under a logarithmic loss constraint, vector Gaussian quantization codebooks with time-sharing are optimal. In the case of quadratic distortion constraint, however, a characterization of the rate-distortion region is still to be found in general, and it is not known yet whether vector Gaussian quantization codebooks (with or without time-sharing) are optimal, except in few special cases such as that of scalar Gaussian sources or the case of only one agent, i.e., the remote vector Gaussian Wyner-Ziv problem whose rate-distortion region is found in [35]. In [35], Tian and Chen also found the rate-distortion region of the remote vector Gaussian Wyner-Ziv problem under logarithmic loss, which they showed achievable using Gaussian quantization codebooks that are different from those (also Gaussian) that are optimal in the case of quadratic distortion. As we already mentioned, our result of Theorem 2 generalizes that of [35] to the case of an arbitrary number of agents.  $\blacksquare$

**Remark 4.** One may wonder whether giving the decoder side information  $\mathbf{Y}_0$  to the encoders is beneficial. Similar to the well known result in Wyner-Ziv source coding of scalar Gaussian sources, our result of Theorem 2 shows that encoder side information does not help.  $\blacksquare$

#### IV. APPLICATIONS

In this section, we show that application of the result of Theorem 2 allows us to find the complete solutions of two related problems: a quadratic vector Gaussian CEO problem with determinant constraint and the vector Gaussian distributed Information Bottleneck problem. For the case of discrete data, we provide an example application to distributed pattern classification.

##### A. Quadratic Vector Gaussian CEO Problem with Determinant Constraint

We now turn to the case in which the distortion is measured under quadratic loss. In this case, the mean square error matrix is defined by

$$\mathbf{D}^{(n)} := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\mathbf{X}_i - \hat{\mathbf{X}}_i)(\mathbf{X}_i - \hat{\mathbf{X}}_i)^\dagger]. \quad (12)$$

Under a (general) error constraint of the form

$$\mathbf{D}^{(n)} \preceq \mathbf{D}, \quad (13)$$

where  $\mathbf{D}$  designates here a prescribed positive definite error matrix, a complete solution is still to be found in general. In what follows, we replace the constraint (13) with one on the *determinant* of the error matrix  $\mathbf{D}^{(n)}$ , i.e.,

$$|\mathbf{D}^{(n)}| \leq D, \quad (14)$$

( $D$  is a scalar here). We note that since the error matrix  $\mathbf{D}^{(n)}$  is minimized by choosing the decoding as

$$\hat{\mathbf{X}}_i = \mathbb{E}[\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n], \quad (15)$$

where  $\{\check{\phi}_k^{(n)}\}_{k=1}^K$  denote the encoding functions, without loss of generality we can write (12) as

$$\mathbf{D}^{(n)} = \frac{1}{n} \sum_{i=1}^n \text{mmse}(\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n). \quad (16)$$

**Definition 5.** A rate-distortion tuple  $(R_1, \dots, R_K, D)$  is achievable for the quadratic vector Gaussian CEO problem with determinant constraint if there exist a blocklength  $n$ ,  $K$  encoding functions  $\{\check{\phi}_k^{(n)}\}_{k=1}^K$  such that

$$\begin{aligned} R_k &\geq \frac{1}{n} \log M_k^{(n)}, \quad \text{for } k = 1, \dots, K, \\ D &\geq \left| \frac{1}{n} \sum_{i=1}^n \text{mmse}(\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \right|. \end{aligned}$$

The rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^{\det}$  is defined as the closure of all non-negative tuples  $(R_1, \dots, R_K, D)$  that are achievable. ■

The following theorem characterizes the rate-distortion region of the quadratic vector Gaussian CEO problem with determinant constraint.

**Theorem 3.** The rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^{\det}$  of the quadratic vector Gaussian CEO problem with determinant constraint is given by the set of all non-negative tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$\log \frac{1}{D} \leq \sum_{k \in \mathcal{S}} R_k + \log |\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k| + \log \left| \mathbf{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{H}_{\bar{\mathcal{S}}}^{\dagger} \mathbf{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1} (\mathbf{I} - \mathbf{\Lambda}_{\mathcal{S}} \mathbf{\Sigma}_{\mathbf{n}_{\mathcal{S}}}^{-1}) \mathbf{H}_{\mathcal{S}} \right|,$$

for matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ , where  $\bar{\mathcal{S}} = \{0\} \cup \mathcal{S}^c$  and  $\mathbf{\Lambda}_{\bar{\mathcal{S}}}$  is as defined by (11).

*Proof.* The proof of Theorem 3 is given in Appendix III. □

**Remark 5.** It is believed that the approach of this section, which connects the quadratic vector Gaussian CEO problem to that under logarithmic loss, can also be exploited to possibly infer other new results on the quadratic vector Gaussian CEO problem. Alternatively, it can also be used to derive new converses on the quadratic vector

*Gaussian CEO problem.* For example, in the case of scalar sources Theorem 3 and Lemma 8 readily provide an alternate converse proof to those of [3], [4] for this model. Similar connections were made in [47], [48] where it was observed that the results of [10] can be used to recover known results on the scalar Gaussian CEO problem (such as the sum rate-distortion region of [49]) and the scalar Gaussian two-encoder distributed source coding problem. We also point out that similar information constraints have been applied to log-determinant reproduction constraints previously in [50]. ■

### B. Distributed Vector Gaussian Information Bottleneck

Consider now the vector Gaussian CEO problem with side information of Section III, and let the logarithmic loss distortion constraint be replaced by the mutual information constraint

$$I\left(\mathbf{X}^n; \psi^{(n)}\left(\phi_1^{(n)}(\mathbf{Y}_1^n), \dots, \phi_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n\right)\right) \geq n\Delta. \quad (17)$$

In this case, the region of optimal tuples  $(R_1, \dots, R_K, \Delta)$  generalizes the *Gaussian Information Bottleneck Function* of [26]–[28] to the setting in which the decoder observes correlated side information  $\mathbf{Y}_0$  and the inference is done in a distributed manner by  $K$  learners. This region can be obtained readily from Theorem 2 by substituting therein  $\Delta := h(\mathbf{X}) - D$ . The following corollary states the result.

**Corollary 2.** *For the problem of distributed Gaussian Information Bottleneck with side information at the predictor, the complexity-relevance region is given by the union of all non-negative tuples  $(R_1, \dots, R_K, \Delta)$  that satisfy, for every  $\mathcal{S} \subseteq \mathcal{K}$ ,*

$$\Delta \leq \sum_{k \in \mathcal{S}} (R_k + \log |\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k|) + \log |\mathbf{I} + \mathbf{\Sigma}_x \mathbf{H}_{\mathcal{S}}^\dagger \mathbf{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1} (\mathbf{I} - \mathbf{\Lambda}_{\bar{\mathcal{S}}} \mathbf{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1}) \mathbf{H}_{\bar{\mathcal{S}}}|,$$

for matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ , where  $\bar{\mathcal{S}} = \{0\} \cup \mathcal{S}^c$  and  $\mathbf{\Lambda}_{\bar{\mathcal{S}}}$  is given by (11). ■

In particular, if  $K = 1$  and  $\mathbf{Y}_0 = \emptyset$ , with the substitutions  $\mathbf{Y} := \mathbf{Y}_1$ ,  $R := R_1$ ,  $\mathbf{H} := \mathbf{H}_1$ ,  $\mathbf{\Sigma} := \mathbf{\Sigma}_1$ , and  $\mathbf{\Omega}_1 := \mathbf{\Omega}$ , the region of Corollary 2 reduces to the set of pairs  $(R, \Delta)$  that satisfy

$$\Delta \leq \log |\mathbf{I} + \mathbf{\Sigma}_x \mathbf{H}^\dagger \mathbf{\Omega} \mathbf{H}| \quad (18a)$$

$$\Delta \leq R + \log |\mathbf{I} - \mathbf{\Omega} \mathbf{\Sigma}|, \quad (18b)$$

for some matrix  $\mathbf{\Omega}$  such that  $\mathbf{0} \preceq \mathbf{\Omega} \preceq \mathbf{\Sigma}^{-1}$ .

Expression (18) is known as the *Gaussian Information Bottleneck Function* [26]–[28], which is the solution of the Information Bottleneck method of [14] in the case of jointly Gaussian variables. More precisely, using the terminology of [14], the inequalities (18) describe the optimal trade-off between the complexity (or rate)  $R$  and the relevance (or accuracy)  $\Delta$ . The concept of Information Bottleneck was found useful in various learning applications, such as for data clustering [51], feature selection [52] and others, including in distributed settings [53], [54].

Furthermore, if in (3) and (4) the noises are independent among them and from  $\mathbf{N}_0$ , the relevance-complexity region of Corollary 2 reduces to the union of all non-negative tuples  $(R_1, \dots, R_K, \Delta)$  that satisfy, for every  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$\Delta \leq \sum_{k \in \mathcal{S}} (R_k + \log |\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k|) + \log |\mathbf{I} + \mathbf{\Sigma}_x (\mathbf{H}_0^\dagger \mathbf{\Sigma}_0^{-1} \mathbf{H}_0 + \sum_{k \in \mathcal{S}^c} \mathbf{H}_k^\dagger \mathbf{\Omega}_k \mathbf{H}_k)|,$$

for some matrices  $\{\mathbf{\Omega}_k\}_{k=1}^K$  such that  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ .

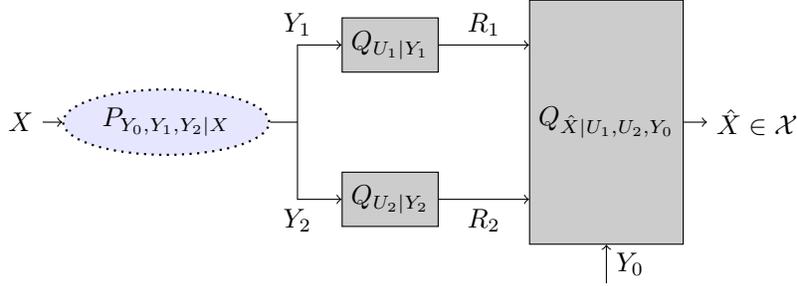


Fig. 2: An example of distributed pattern classification.

### C. Distributed Pattern Classification

Consider the problem of distributed pattern classification shown in Figure 2. In this example, the decoder is a predictor whose role is to guess the unknown class  $X \in \mathcal{X}$  of a measurable pair  $(Y_1, Y_2) \in \mathcal{Y}_1 \times \mathcal{Y}_2$  on the basis of inputs from two learners as well as its own observation about the target class, in the form of some correlated  $Y_0 \in \mathcal{Y}_0$ . It is assumed that  $Y_1 \oslash (X, Y_0) \oslash Y_2$ . The first learner produces its input based only on  $Y_1 \in \mathcal{Y}_1$ ; and the second learner produces its input based only on  $Y_2 \in \mathcal{Y}_2$ . For the sake of a smaller *generalization gap*<sup>1</sup>, the inputs of the learners are restricted to have description lengths that are no more than  $R_1$  and  $R_2$  bits per sample, respectively. Let  $Q_{U_1 | Y_1} : \mathcal{Y}_1 \rightarrow \mathcal{P}(\mathcal{U}_1)$  and  $Q_{U_2 | Y_2} : \mathcal{Y}_2 \rightarrow \mathcal{P}(\mathcal{U}_2)$  be two (stochastic) such learners. Also, let  $Q_{\hat{X} | U_1, U_2, Y_0} : \mathcal{U}_1 \times \mathcal{U}_2 \times \mathcal{Y}_0 \rightarrow \mathcal{P}(\mathcal{X})$  be a soft-decoder or predictor that maps the pair of representations  $(U_1, U_2)$  and  $Y_0$  to a probability distribution on the label space  $\mathcal{X}$ . The pair of learners and predictor induce a classifier

$$\begin{aligned} Q_{\hat{X} | Y_0, Y_1, Y_2}(x | y_0, y_1, y_2) &= \sum_{u_1 \in \mathcal{U}_1} Q_{U_1 | Y_1}(u_1 | y_1) \sum_{u_2 \in \mathcal{U}_2} Q_{U_2 | Y_2}(u_2 | y_2) Q_{\hat{X} | U_1, U_2, Y_0}(x | u_1, u_2, y_0) \\ &= \mathbb{E}_{Q_{U_1 | Y_1}} \mathbb{E}_{Q_{U_2 | Y_2}} [Q_{\hat{X} | U_1, U_2, Y_0}(x | U_1, U_2, y_0)], \end{aligned} \quad (19)$$

whose probability of classification error is defined as

$$P_{\mathcal{E}}(Q_{\hat{X} | Y_0, Y_1, Y_2}) = 1 - \mathbb{E}_{P_{X, Y_0, Y_1, Y_2}} [Q_{\hat{X} | Y_0, Y_1, Y_2}(X | Y_0, Y_1, Y_2)]. \quad (20)$$

<sup>1</sup>The generalization gap, defined as the difference between the empirical risk (average risk over a finite training sample) and the population risk (average risk over the true joint distribution), can be upper bounded using the mutual information between the learner's inputs and outputs, see, e.g., [55], [56] and the recent [57], which provides a fundamental justification of the use of the *minimum description length* (MDL) constraint on the learners mappings as a regularizer term.

Let  $\mathcal{RD}_{\text{CEO}}^*$  be the rate-distortion region of the associated two-encoder DM CEO problem with side information as given by Theorem 1. The following proposition shows that there exists a classifier  $Q_{\hat{X}|Y_0, Y_1, Y_2}^*$  for which the probability of misclassification can be upper bounded in terms of the minimal average logarithmic loss distortion that is achievable for the rate pair  $(R_1, R_2)$  in  $\mathcal{RD}_{\text{CEO}}^*$ .

**Proposition 2.** *For the problem of distributed pattern classification of Figure 2, there exists a classifier  $Q_{\hat{X}|Y_0, Y_1, Y_2}^*$  for which the probability of classification error satisfies*

$$P_{\mathcal{E}}(Q_{\hat{X}|Y_0, Y_1, Y_2}^*) \leq 1 - \exp(-\inf\{D : (R_1, R_2, D) \in \mathcal{RD}_{\text{CEO}}^*\}) ,$$

where  $\mathcal{RD}_{\text{CEO}}^*$  is the rate-distortion region of the associated two-encoder DM CEO problem with side information as given by Theorem 1.

*Proof.* The proof of Proposition 2 is given in Appendix IV. □

To make the above example more concrete, consider the following scenario where  $Y_0$  plays the role of information about the sub-class of the label class  $X \in \{0, 1, 2, 3\}$ . More specifically, let  $S$  be a random variable that is uniformly distributed over  $\{1, 2\}$ . Also, let  $X_1$  and  $X_2$  be two random variables that are independent between them and from  $S$ , distributed uniformly over  $\{1, 3\}$  and  $\{0, 2\}$  respectively. The state  $S$  acts as a random switch that connects  $X_1$  or  $X_2$  to  $X$ , i.e.,

$$X = X_S . \tag{21}$$

That is, if  $S = 1$  then  $X = X_1$ , and if  $S = 2$  then  $X = X_2$ . Thus, the value of  $S$  indicates whether  $X$  is odd- or even-valued (i.e., the sub-class of  $X$ ). Also, let

$$Y_0 = S \tag{22a}$$

$$Y_1 = X_S \oplus Z_1 \tag{22b}$$

$$Y_2 = X_S \oplus Z_2 , \tag{22c}$$

where  $Z_1$  and  $Z_2$  are Bernoulli- $(p)$  random variables,  $p \in (0, 1)$ , that are independent between them, and from  $(S, X_1, X_2)$ , and the addition is modulo 4. For simplification, we let  $R_1 = R_2 = R$ . We numerically approximate the set of  $(R, D)$  pairs such that  $(R, R, D)$  is in the rate-distortion region  $\mathcal{RD}_{\text{CEO}}^*$  corresponding to the CEO network of this example. The algorithm that we use for the computation will be described in detail in Section V-A. The lower convex envelope of these  $(R, D)$  pairs is plotted in Figure 3a for  $p \in \{0.01, 0.1, 0.25, 0.5\}$ . Continuing our example, we also compute the upper bound on the probability of classification error according to Proposition 2. The result is given in Figure 3b. Observe that if  $Y_1$  and  $Y_2$  are high-quality estimates of  $X$  (e.g.,  $p = 0.01$ ), then a small increase in the *complexity*  $R$  results in a large relative improvement of the (bound on) the probability of classification error. On the other hand, if  $Y_1$  and  $Y_2$  are low-quality estimates of  $X$  (e.g.,  $p = 0.25$ ) then we

require a large increase of  $R$  in order to obtain an appreciable reduction in the error probability. Recalling that larger  $R$  implies lesser generalization capability [55]–[57], these numerical results are consistent with the fact that classifiers should strike a good balance between accuracy and their ability to generalize well to unseen data. Figure 3c quantifies the value of side information  $S$  given to both learners and predictor, none of them, or only the predictor, for  $p = 0.25$ .

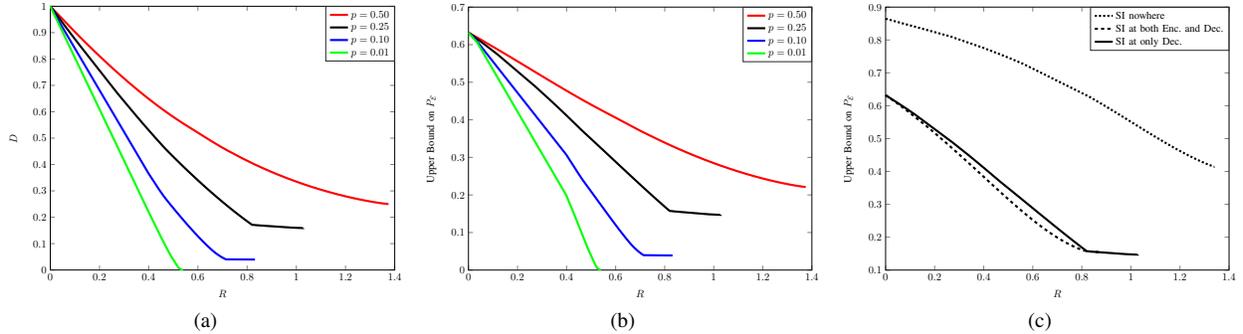


Fig. 3: Illustration of the bound on the probability of classification error of Proposition 2 for the example described by (21) and (22). (a) Distortion-rate function of the network of Figure 2 computed for  $p \in \{0.01, 0.1, 0.25, 0.5\}$ . (b) Upper bound on the probability of classification error computed according to Proposition 2. (c) Effect of side information (SI)  $Y_0$  when given to both learners and the predictor, only the predictor or none of them.

## V. BLAHUT-ARIMOTO TYPE ALGORITHMS

In this section, we develop iterative algorithms that allow to compute the rate-distortion regions of the DM and vector Gaussian CEO problems numerically. We illustrate the efficiency of our algorithms through some numerical examples.

### A. Discrete Case

Here we develop a BA-type algorithm that allows to compute the convex region  $\mathcal{RD}_{\text{CEO}}^*$  for general discrete memoryless sources. To develop the algorithm, we use the Berger-Tung form of the region given in Definition 4 for  $K = 2$ . The outline of the proposed method is as follows. First, we rewrite the rate-distortion region  $\mathcal{RD}_{\text{CEO}}^*$  in terms of the union of two simpler regions in Proposition 3. The tuples lying on the boundary of each region are parametrically given in Theorem 4. Then, the boundary points of each simpler region are computed numerically via an alternating minimization method derived in Section V-A2 and detailed in Algorithm 1. Finally, the original rate-distortion region is obtained as the convex hull of the union of the tuples obtained for the two simple regions.

#### 1) Equivalent Parametrization:

Define the two regions  $\mathcal{RD}_{\text{CEO}}^k$ ,  $k = 1, 2$ , as

$$\mathcal{RD}_{\text{CEO}}^k = \{(R_1, R_2, D) : D \geq D_{\text{CEO}}^k(R_1, R_2)\}, \quad (23)$$

with

$$D_{\text{CEO}}^k(R_1, R_2) := \min H(X|U_1, U_2, Y_0) \quad (24)$$

$$\text{s.t. } R_k \geq I(Y_k; U_k|U_{\bar{k}}, Y_0) \quad \text{and} \quad R_{\bar{k}} \geq I(X_{\bar{k}}; U_{\bar{k}}|Y_0),$$

and the minimization is over set of joint measures  $P_{U_1, U_2, X, Y_0, Y_1, Y_2}$  that satisfy  $U_1 \circlearrowleft Y_1 \circlearrowleft (X, Y_0) \circlearrowleft Y_2 \circlearrowleft U_2$ . (We define  $\bar{k} := k \pmod{2} + 1$  for  $k = 1, 2$ .)

As stated in the following proposition, the region  $\mathcal{RD}_{\text{CEO}}^*$  of Theorem 1 coincides with the convex hull of the union of the two regions  $\mathcal{RD}_{\text{CEO}}^1$  and  $\mathcal{RD}_{\text{CEO}}^2$ .

**Proposition 3.** *The region  $\mathcal{RD}_{\text{CEO}}^*$  is given by*

$$\mathcal{RD}_{\text{CEO}}^* = \text{conv}(\mathcal{RD}_{\text{CEO}}^1 \cup \mathcal{RD}_{\text{CEO}}^2). \quad (25)$$

*Proof.* An outline of the proof is as follows. Let  $P_{U_1, U_2, X, Y_0, Y_1, Y_2}$  and  $P_Q$  be such that  $(R_1, R_2, D) \in \mathcal{RD}_{\text{CEO}}^*$ . The polytope defined by the rate constraints (9), denoted by  $\mathcal{V}$ , forms a contra-polymatroid with  $2!$  extreme points (vertices) [10], [58]. Given a permutation  $\pi$  on  $\{1, 2\}$ , the tuple

$$\tilde{R}_{\pi(1)} = I(Y_{\pi(1)}; U_{\pi(1)}|Y_0), \quad \tilde{R}_{\pi(2)} = I(Y_{\pi(2)}; U_{\pi(2)}|U_{\pi(1)}, Y_0),$$

defines an extreme point of  $\mathcal{V}$  for each permutation. As shown in [10], for every extreme point  $(\tilde{R}_1, \tilde{R}_2)$  of  $\mathcal{V}$ , the point  $(\tilde{R}_1, \tilde{R}_2, D)$  is achieved by time-sharing two successive Wyner-Ziv (WZ) strategies. The set of achievable tuples with such successive WZ scheme is characterized by the convex hull of  $\mathcal{RD}_{\text{CEO}}^{\pi(1)}$ . Convexifying the union of both regions as in (25), we obtain the full rate-distortion region  $\mathcal{RD}_{\text{CEO}}^*$ .  $\square$

The main advantage of Proposition 3 is that it reduces the computation of region  $\mathcal{RD}_{\text{CEO}}^*$  to the computation of the two regions  $\mathcal{RD}_{\text{CEO}}^k$ ,  $k = 1, 2$ , whose boundary can be efficiently parametrized, leading to an efficient computational method. In what follows, we concentrate on  $\mathcal{RD}_{\text{CEO}}^1$ . The computation of  $\mathcal{RD}_{\text{CEO}}^2$  follows similarly, and is omitted for brevity. Next theorem provides a parametrization of the boundary tuples of the region  $\mathcal{RD}_{\text{CEO}}^1$  in terms, each of them, of an optimization problem over the pmfs  $\mathbf{P} := \{P_{U_1|Y_1}, P_{U_2|Y_2}\}$ .

**Theorem 4.** *For each  $\mathbf{s} := [s_1, s_2]$ ,  $s_1 > 0$ ,  $s_2 > 0$ , define a tuple  $(R_{1,\mathbf{s}}, R_{2,\mathbf{s}}, D_{\mathbf{s}})$  parametrically given by*

$$D_{\mathbf{s}} = -s_1 R_{1,\mathbf{s}} - s_2 R_{2,\mathbf{s}} + \min_{\mathbf{P}} F_{\mathbf{s}}(\mathbf{P}), \quad (26)$$

$$R_{1,\mathbf{s}} = I(Y_1; U_1^*|U_2^*, Y_0), \quad R_{2,\mathbf{s}} = I(Y_2; U_2^*|Y_0), \quad (27)$$

where  $F_{\mathbf{s}}(\mathbf{P})$  is given as follows

$$F_{\mathbf{s}}(\mathbf{P}) := H(X|U_1, U_2, Y_0) + s_1 I(Y_1; U_1|U_2, Y_0) + s_2 I(Y_2; U_2|Y_0), \quad (28)$$

and;  $\mathbf{P}^*$  are the conditional pmfs yielding the minimum in (26) and  $U_1^*, U_2^*$  are the auxiliary variables induced by  $\mathbf{P}^*$ . Then, we have:

- 1) Each value of  $\mathbf{s}$  leads to a tuple  $(R_{1,\mathbf{s}}, R_{2,\mathbf{s}}, D_{\mathbf{s}})$  on the distortion-rate curve  $D_{\mathbf{s}} = D_{\text{CEO}}^1(R_{1,\mathbf{s}}, R_{2,\mathbf{s}})$ .
- 2) For every point on the distortion-rate curve, there is an  $\mathbf{s}$  for which (26) and (27) hold.

*Proof.* Suppose that  $\mathbf{P}^*$  yields the minimum in (26). For this  $\mathbf{P}$ , we have  $I(Y_1; U_1|U_2, Y_0) = R_{1,\mathbf{s}}$  and  $I(Y_2; U_2|Y_0) = R_{2,\mathbf{s}}$ . Then, we have

$$\begin{aligned} D_{\mathbf{s}} &= -s_1 R_{1,\mathbf{s}} - s_2 R_{2,\mathbf{s}} + F_{\mathbf{s}}(\mathbf{P}^*) \\ &= -s_1 R_{1,\mathbf{s}} - s_2 R_{2,\mathbf{s}} + [H(X|U_1^*, U_2^*, Y_0) + s_1 R_{1,\mathbf{s}} + s_2 R_{2,\mathbf{s}}] \\ &= H(X|U_1^*, U_2^*, Y_0) \geq D_{\text{CEO}}^1(R_{1,\mathbf{s}}, R_{2,\mathbf{s}}). \end{aligned} \quad (29)$$

Conversely, if  $\mathbf{P}^*$  is the solution to the minimization in (24), then  $I(Y_1; U_1^*|U_2^*, Y_0) \leq R_1$  and  $I(Y_2; U_2^*|Y_0) \leq R_2$  and for any  $\mathbf{s}$ ,

$$\begin{aligned} D_{\text{CEO}}^1(R_1, R_2) &= H(X|U_1^*, U_2^*, Y_0) \\ &\geq H(X|U_1^*, U_2^*, Y_0) + s_1(I(Y_1; U_1^*|U_2^*, Y_0) - R_1) + s_2(I(Y_2; U_2^*|Y_0) - R_2) \\ &= D_{\mathbf{s}} + s_1(R_{1,\mathbf{s}} - R_1) + s_2(R_{2,\mathbf{s}} - R_2). \end{aligned}$$

Given  $\mathbf{s}$ , and hence  $(R_{1,\mathbf{s}}, R_{2,\mathbf{s}}, D_{\mathbf{s}})$ , letting  $(R_1, R_2) = (R_{1,\mathbf{s}}, R_{2,\mathbf{s}})$  yields  $D_{\text{CEO}}^1(R_{1,\mathbf{s}}, R_{2,\mathbf{s}}) \geq D_{\mathbf{s}}$ , which proves, together with (29), statement 1) and 2).  $\square$

Next, we show that it is sufficient to run the algorithm for  $s_1 \in (0, 1]$ .

**Lemma 1.** *The range of the parameter  $s_1$  can be restricted to  $(0, 1]$ .*

*Proof.* Let  $F^* = \min_{\mathbf{P}} F_{\mathbf{s}}(\mathbf{P})$ . If we set  $U_1 = \emptyset$ , then we have the relation  $F^* \leq H(X|U_2, Y_0) + s_2 I(Y_2; U_2|Y_0)$ . For  $s_1 > 1$ , we have

$$F_{\mathbf{s}}(\mathbf{P}) \stackrel{(a)}{\geq} (1 - s_1)H(X|U_1, U_2, Y_0) + s_1 H(X|U_2, Y_0) + s_2 I(Y_2; U_2|Y_0) \stackrel{(b)}{\geq} H(X|U_2, Y_0) + s_2 I(Y_2; U_2|Y_0),$$

where (a) follows since mutual information is always positive, i.e.,  $I(Y_1; U_1|X, Y_0) \geq 0$ ; (b) holds since conditioning reduces entropy and  $1 - s_1 < 0$ . Then  $F^* = H(X|U_2, Y_0) + s_2 I(Y_2; U_2|Y_0)$  for  $s_1 > 1$ . Hence we can restrict the range of  $s_1$  to  $s_1 \in (0, 1]$ .  $\square$

2) *Computation of  $\mathcal{RD}_{\text{CEO}}^1$ :*

In this section, we derive an algorithm to solve (26) for a given parameter value  $\mathbf{s}$ . To that end, we define a variational bound on  $F_{\mathbf{s}}(\mathbf{P})$ , and optimize it instead of (26). Let  $\mathbf{Q}$  be a set of some auxiliary pmfs defined as

$$\mathbf{Q} := \{Q_{U_1}, Q_{U_2}, Q_{X|U_1, U_2, Y_0}, Q_{X|U_1, Y_0}, Q_{X|U_2, Y_0}, Q_{Y_0|U_1}, Q_{Y_0|U_2}\}. \quad (30)$$

In the following we define the variational cost function  $F_s(\mathbf{P}, \mathbf{Q})$

$$\begin{aligned} F_s(\mathbf{P}, \mathbf{Q}) := & -s_1 H(X|Y_0) - (s_1 + s_2) H(Y_0) \\ & + \mathbb{E}_{P_{X,Y_0,Y_1,Y_2}} \left[ (1-s_1) \mathbb{E}_{P_{U_1|Y_1}} \mathbb{E}_{P_{U_2|Y_2}} [-\log Q_{X|U_1,U_2,Y_0}] + s_1 \mathbb{E}_{P_{U_1|Y_1}} [-\log Q_{X|U_1,Y_0}] \right. \\ & + s_1 \mathbb{E}_{P_{U_2|Y_2}} [-\log Q_{X|U_2,Y_0}] + s_1 D_{\text{KL}}(P_{U_1|Y_1} \| Q_{U_1}) + s_2 D_{\text{KL}}(P_{U_2|Y_2} \| Q_{U_2}) \\ & \left. + s_1 \mathbb{E}_{P_{U_1|Y_1}} [-\log Q_{Y_0|U_1}] + s_2 \mathbb{E}_{P_{U_2|Y_2}} [-\log Q_{Y_0|U_2}] \right]. \end{aligned} \quad (31)$$

The following lemma states that  $F_s(\mathbf{P}, \mathbf{Q})$  is an upper bound on  $F_s(\mathbf{P})$  for all distributions  $\mathbf{Q}$ .

**Lemma 2.** *For fixed  $\mathbf{P}$ , we have*

$$F_s(\mathbf{P}, \mathbf{Q}) \geq F_s(\mathbf{P}), \quad \text{for all } \mathbf{Q}.$$

*In addition, there exists a  $\mathbf{Q}$  that achieves the minimum  $\min_{\mathbf{Q}} F_s(\mathbf{P}, \mathbf{Q}) = F_s(\mathbf{P})$ , given by*

$$\begin{aligned} Q_{U_k} &= P_{U_k}, \quad Q_{X|U_k,Y_0} = P_{X|U_k,Y_0}, \quad Q_{Y_0|U_k} = P_{Y_0|U_k}, \quad \text{for } k = 1, 2, \\ Q_{X|U_1,U_2,Y_0} &= P_{X|U_1,U_2,Y_0}. \end{aligned} \quad (32)$$

*Proof.* The proof of Lemma 2 is given in Appendix V. □

Using the lemma above, the minimization in (26) can be written in terms of the variational cost function as follows

$$\min_{\mathbf{P}} F_s(\mathbf{P}) = \min_{\mathbf{P}} \min_{\mathbf{Q}} F_s(\mathbf{P}, \mathbf{Q}). \quad (33)$$

Motivated by the BA algorithm [24], [25], we propose an alternate optimization procedure over the set of pmfs  $\mathbf{P}$  and  $\mathbf{Q}$  as stated in Algorithm 1. The main idea is that at iteration  $t$ , for fixed  $\mathbf{P}^{(t-1)}$  the optimal  $\mathbf{Q}^{(t)}$  minimizing  $F_s(\mathbf{P}, \mathbf{Q})$  can be found analytically; next, for given  $\mathbf{Q}^{(t)}$  the optimal  $\mathbf{P}^{(t)}$  that minimizes  $F_s(\mathbf{P}, \mathbf{Q})$  has also a closed form. So, starting with a random initialization  $\mathbf{P}^{(0)}$ , the algorithm iterates over distributions  $\mathbf{Q}$  and  $\mathbf{P}$  minimizing  $F_s(\mathbf{P}, \mathbf{Q})$  until the convergence, as stated below

$$\mathbf{P}^{(0)} \rightarrow \mathbf{Q}^{(1)} \rightarrow \mathbf{P}^{(1)} \rightarrow \dots \rightarrow \mathbf{P}^{(t-1)} \rightarrow \mathbf{Q}^{(t)} \rightarrow \dots \rightarrow \mathbf{P}^* \rightarrow \mathbf{Q}^*.$$

At each iteration, the optimal values of  $\mathbf{P}$  and  $\mathbf{Q}$  are found by solving a convex optimization problems. We have the following lemma.

**Lemma 3.**  *$F_s(\mathbf{P}, \mathbf{Q})$  is convex in  $\mathbf{P}$  and convex in  $\mathbf{Q}$ .*

*Proof.* The proof of Lemma 3 follows from the log-sum inequality. □

For fixed  $\mathbf{P}^{(t-1)}$ , the optimal  $\mathbf{Q}^{(t)}$  minimizing the variational bound in (31) can be found from Lemma 2 and given by (32). For fixed  $\mathbf{Q}^{(t)}$ , the optimal  $\mathbf{P}^{(t)}$  minimizing (31) can be found by using the next lemma.

**Algorithm 1** BA-type algorithm to compute  $\mathcal{RD}_{\text{CEO}}^1$ 

- 1: **input:** pmf  $P_{X,Y_0,Y_1,Y_2}$ , parameters  $1 \geq s_1 > 0$ ,  $s_2 > 0$ .
- 2: **output:** Optimal  $P_{U_1|Y_1}^*$ ,  $P_{U_2|Y_2}^*$ ; triple  $(R_{1,s}, R_{2,s}, D_s)$ .
- 3: **initialization** Set  $t = 0$ . Set  $\mathbf{P}^{(0)}$  randomly.
- 4: **repeat**
- 5:     Update the following pmfs for  $k = 1, 2$

$$\begin{aligned} p^{(t+1)}(u_k) &= \sum_{y_k} p^{(t)}(u_k|y_k)p(y_k), \\ p^{(t+1)}(u_k|y_0) &= \sum_{y_k} p^{(t)}(u_k|y_k)p(y_k|y_0), \\ p^{(t+1)}(u_k|x, y_0) &= \sum_{y_k} p^{(t)}(u_k|y_k)p(y_k|x, y_0), \\ p^{(t+1)}(x|u_1, u_2, y_0) &= \frac{p^{(t+1)}(u_1|x, y_0)p^{(t+1)}(u_2|x, y_0)p(x, y_0)}{\sum_x p^{(t+1)}(u_1|x, y_0)p^{(t+1)}(u_2|x, y_0)p(x, y_0)}. \end{aligned}$$

- 6:     Update  $\mathbf{Q}^{(t+1)}$  by using (32).
- 7:     Update  $\mathbf{P}^{(t+1)}$  by using (34).
- 8:      $t \leftarrow t + 1$ .
- 9: **until** convergence.

**Lemma 4.** For fixed  $\mathbf{Q}$ , there exists a  $\mathbf{P}$  that achieves the minimum  $\min_{\mathbf{P}} F_s(\mathbf{P}, \mathbf{Q})$ , where  $P_{U_k|Y_k}$  is given by

$$p(u_k|y_k) = q(u_k) \frac{\exp[-\psi_k(u_k, y_k)]}{\sum_{u_k} q(u_k) \exp[-\psi_k(u_k, y_k)]}, \quad \text{for } k = 1, 2, \quad (34)$$

where  $\psi_k(u_k, y_k)$ ,  $k = 1, 2$ , are defined as follows

$$\psi_k(u_k, y_k) := \frac{1 - s_1}{s_k} \mathbb{E}_{U_{\bar{k}}, Y_0|y_k} D(P_{X|y_k, U_{\bar{k}}, Y_0} \| Q_{X|u_k, U_{\bar{k}}, Y_0}) + \frac{s_1}{s_k} \mathbb{E}_{Y_0|y_k} D(P_{X|y_k, Y_0} \| Q_{X|u_k, Y_0}) + D(P_{Y_0|y_k} \| Q_{Y_0|u_k}). \quad (35)$$

*Proof.* The proof of Lemma 4 is given in Appendix VI.  $\square$

At each iteration of Algorithm 1,  $F_s(\mathbf{P}^{(t)}, \mathbf{Q}^{(t)})$  decreases until eventually it converges. However, since  $F_s(\mathbf{P}, \mathbf{Q})$  is convex in each argument but not necessarily jointly convex, Algorithm 1 does not necessarily converge to the global optimum. In particular, next proposition shows that Algorithm 1 converges to a stationary solution of the minimization in (26).

**Proposition 4.** Every limit point of  $\mathbf{P}^{(t)}$  generated by Algorithm 1 converges to a stationary solution of (26).

*Proof.* Algorithm 1 falls into the class of so-called ‘‘Successive Upper-bound Minimization’’ (SUM) algorithms [59], in which  $F_s(\mathbf{P}, \mathbf{Q})$  acts as a globally tight upper bound on  $F_s(\mathbf{P})$ . Let  $\mathbf{Q}^*(\mathbf{P}) := \arg \min_{\mathbf{Q}} F_s(\mathbf{P}, \mathbf{Q})$ . From Lemma 2,  $F_s(\mathbf{P}, \mathbf{Q}^*(\mathbf{P}')) \geq F_s(\mathbf{P}, \mathbf{Q}^*(\mathbf{P})) = F_s(\mathbf{P})$  for  $\mathbf{P}' \neq \mathbf{P}$ . It follows that  $F_s(\mathbf{P})$  and  $F_s(\mathbf{P}, \mathbf{Q}^*(\mathbf{P}'))$  satisfy [59, Proposition 1] and thus  $F_s(\mathbf{P}, \mathbf{Q}^*(\mathbf{P}'))$  satisfies (A1)–(A4) in [59]. Convergence to a stationary point of (26) follows from [59, Theorem 1].  $\square$

**Remark 6.** Algorithm 1 generates a sequence that is non-increasing. Since this sequence is lower bounded, convergence to a stationary point is guaranteed. This per-se, however, does not necessarily imply that such a

point is a stationary solution of the original problem described by (26). Instead, this is guaranteed here by showing that the Algorithm 1 is of SUM-type with the function  $F_s(\mathbf{P}, \mathbf{Q})$  satisfying the necessary conditions [59, (A1)–(A4)]. ■

---

**Algorithm 2** BA-type algorithm for the Gaussian vector CEO

---

- 1: **input:** Covariance  $\Sigma_{(x, y_0, y_1, y_2)}$ , parameters  $s_1 \geq s_2 > 0$ ,  $s_2 > 0$ .
- 2: **output:** Optimal pairs  $(\mathbf{A}_k^*, \Sigma_{z_k^*})$ ,  $k = 1, 2$ .
- 3: **initialization** Set  $t = 0$ . Set randomly  $\mathbf{A}_k^0$  and  $\Sigma_{z_k^0} \succeq 0$  for  $k = 1, 2$ .
- 4: **repeat**
- 5:     For  $k = 1, 2$ , update the following

$$\Sigma_{\mathbf{u}_k^t} = \mathbf{A}_k^t \Sigma_{y_k} \mathbf{A}_k^{t\dagger} + \Sigma_{z_k^t},$$

$$\Sigma_{\mathbf{u}_k^t | (x, y)} = \mathbf{A}_k^t \Sigma_k \mathbf{A}_k^{t\dagger} + \Sigma_{z_k^t},$$

and update  $\Sigma_{\mathbf{u}_k^t | (\mathbf{u}_k^t, y)}$ ,  $\Sigma_{\mathbf{u}_2^t | y}$  and  $\Sigma_{y_k^t | (\mathbf{u}_k^t, y)}$  from their definitions by using the following

$$\Sigma_{\mathbf{u}_1^t, \mathbf{u}_2^t} = \mathbf{A}_1^t \mathbf{H}_1 \Sigma_x \mathbf{H}_2^\dagger \mathbf{A}_2^{t\dagger},$$

$$\Sigma_{\mathbf{u}_k^t, y} = \mathbf{A}_k^t \mathbf{H}_k \Sigma_x \mathbf{H}_0^\dagger,$$

$$\Sigma_{y_k, \mathbf{u}_k^t} = \mathbf{H}_k \Sigma_x \mathbf{H}_k^\dagger \mathbf{A}_k^{t\dagger}.$$

- 6:     Compute  $\Sigma_{z_k^{t+1}}$  as in (38a) for  $k = 1, 2$ .
  - 7:     Compute  $\mathbf{A}_k^{t+1}$  as (38b) for  $k = 1, 2$ .
  - 8:      $t \leftarrow t + 1$ .
  - 9: **until** convergence.
- 

### B. Vector Gaussian Case

Computing the rate-distortion region  $\mathcal{RD}_{\text{VG-CEO}}^*$  of the vector Gaussian CEO problem as given by Theorem 2 is a convex optimization problem on  $\{\Omega_k\}_{k=1}^K$  which can be solved using, e.g., the popular generic optimization tool CVX [60]. Alternatively, the region can be computed using an extension of Algorithm 1 to memoryless Gaussian sources as given in the rest of this section.

For discrete sources with (small) alphabets, the updating rules of  $\mathbf{Q}^{(t+1)}$  and  $\mathbf{P}^{(t+1)}$  of Algorithm 1 are relatively easy computationally. However, they become computationally unfeasible for continuous alphabet sources. Here, we leverage on the optimality of Gaussian test channels as shown by Theorem 2 to restrict the optimization of  $\mathbf{P}$  to Gaussian distributions, which allows to reduce the search of update rules to those of the associated parameters, namely covariance matrices. In particular, we show that if  $P_{\mathbf{U}_k | \mathbf{Y}_k}^{(t)}$ ,  $k = 1, 2$ , is Gaussian and such that

$$\mathbf{U}_k^t = \mathbf{A}_k^t \mathbf{Y}_k + \mathbf{Z}_k^t, \quad (36)$$

where  $\mathbf{Z}_k^t \sim \mathcal{CN}(\mathbf{0}, \Sigma_{z_k^t})$  then  $P_{\mathbf{U}_k | \mathbf{Y}_k}^{(t+1)}$  too is Gaussian, with

$$\mathbf{U}_k^{t+1} = \mathbf{A}_k^{t+1} \mathbf{Y}_k + \mathbf{Z}_k^{t+1}, \quad (37)$$

where  $\mathbf{Z}_k^{t+1} \sim \mathcal{CN}(\mathbf{0}, \Sigma_{\mathbf{z}_k^{t+1}})$  and the parameters  $\mathbf{A}_k^{t+1}$  and  $\Sigma_{\mathbf{z}_k^{t+1}}$  are given by

$$\Sigma_{\mathbf{z}_k^{t+1}} = \left( \frac{1}{s_k} \Sigma_{\mathbf{u}_k^t | (\mathbf{x}, \mathbf{y}_0)}^{-1} - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t | (\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t | \mathbf{y}_0}^{-1} \right)^{-1} \quad (38a)$$

$$\begin{aligned} \mathbf{A}_k^{t+1} = & \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1}{s_k} \Sigma_{\mathbf{u}_k^t | (\mathbf{x}, \mathbf{y}_0)}^{-1} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | (\mathbf{x}, \mathbf{y}_0)} \Sigma_{\mathbf{y}_k}^{-1}) \right) \\ & - \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t | (\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | (\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{\mathbf{y}_k}^{-1}) - \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t | \mathbf{y}_0}^{-1} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | \mathbf{y}_0} \Sigma_{\mathbf{y}_k}^{-1}) \right). \end{aligned} \quad (38b)$$

The updating steps are provided in Algorithm 2. The proof of (38) can be found in Appendix VII.

### C. Numerical Examples

In this section, we discuss two examples, a binary CEO example and a vector Gaussian CEO example.

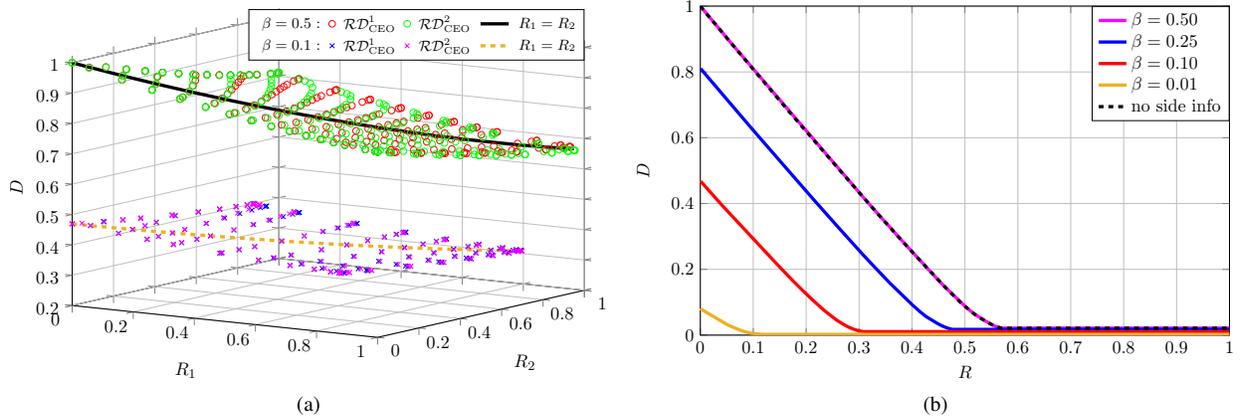


Fig. 4: Rate-distortion region of the binary CEO network of Example 1, computed using Algorithm 1. (a): set of  $(R_1, R_2, D)$  triples such  $(R_1, R_2, D) \in \mathcal{RD}_{\text{CEO}}^1 \cup \mathcal{RD}_{\text{CEO}}^2$ , for  $\alpha_1 = \alpha_2 = 0.25$  and  $\beta \in \{0.1, 0.25\}$ . (b): set of  $(R, D)$  pairs such  $(R, R, D) \in \mathcal{RD}_{\text{CEO}}^1 \cup \mathcal{RD}_{\text{CEO}}^2$ , for  $\alpha_1 = \alpha_2 = 0.01$  and  $\beta \in \{0.01, 0.1, 0.25, 0.5\}$ .

**Example 1.** Consider the following binary CEO problem. A memoryless binary source  $X$ , modeled as a Bernoulli- $(1/2)$  random variable, i.e.,  $X \sim \text{Bern}(1/2)$ , is observed remotely at two agents who communicate with a central unit decoder over error-free rate-limited links of capacity  $R_1$  and  $R_2$ , respectively. The decoder wants to estimate the remote source  $X$  to within some average fidelity level  $D$ , where the distortion is measured under the logarithmic loss criterion. The noisy observation  $Y_1$  at Agent 1 is modeled as the output of a binary symmetric channel (BSC) with crossover probability  $\alpha_1 \in [0, 1]$ , whose input is  $X$ , i.e.,  $Y_1 = X \oplus S_1$  with  $S_1 \sim \text{Bern}(\alpha_1)$ . Similarly, the noisy observation  $Y_2$  at Agent 2 is modeled as the output of a BSC( $\alpha_2$ ) channel,  $\alpha_2 \in [0, 1]$ , whose has input  $X$ , i.e.,  $Y_2 = X \oplus S_2$  with  $S_2 \sim \text{Bern}(\alpha_2)$ . Also, the central unit decoder observes its own side information  $Y_0$  in the

form of the output of a BSC( $\beta$ ) channel,  $\beta \in [0, 1]$ , whose input is  $X$ , i.e.,  $Y_0 = X \oplus S_0$  with  $S_0 \sim \text{Bern}(\beta)$ . It is assumed that the binary noises  $S_0$ ,  $S_1$  and  $S_2$  are independent between them and with the remote source  $X$ .

We use Algorithm 1 to numerically approximate<sup>2</sup> the set of  $(R_1, R_2, D)$  triples such that  $(R_1, R_2, D)$  is in the union of the achievable regions  $\mathcal{RD}_{\text{CEO}}^1$  and  $\mathcal{RD}_{\text{CEO}}^2$  as given by (23). The regions are depicted in Figure 4a for the values  $\alpha_1 = \alpha_2 = 0.25$  and  $\beta \in \{0.1, 0.25\}$ . Note that for both values of  $\beta$ , an approximation of the rate-distortion region  $\mathcal{RD}_{\text{CEO}}$  is easily found as the convex hull of the union of the shown two regions. For simplicity, Figure 4b shows achievable rate-distortion pairs  $(R, D)$  in the case in which the rates of the two encoders are constrained to be at most  $R$  bits per channel use each, i.e.,  $R_1 = R_2 = R$ , higher quality agents' observations  $(Y_1, Y_2)$  corresponding to  $\alpha_1 = \alpha_2 = 0.01$  and  $\beta \in \{0.01, 0.1, 0.25, 0.5\}$ . In this figure, observe that, as expected, smaller values of  $\beta$  correspond to higher quality estimate side information  $Y_0$  at the decoder; and lead to smaller distortion values for given rate  $R$ . The choice  $\beta = 0.5$  corresponds to the case of no or independent side information at decoder; and it is easy to check that the associated  $(R, D)$  curve coincides with the one obtained through exhaustive search in [10, Figure 3]. ■

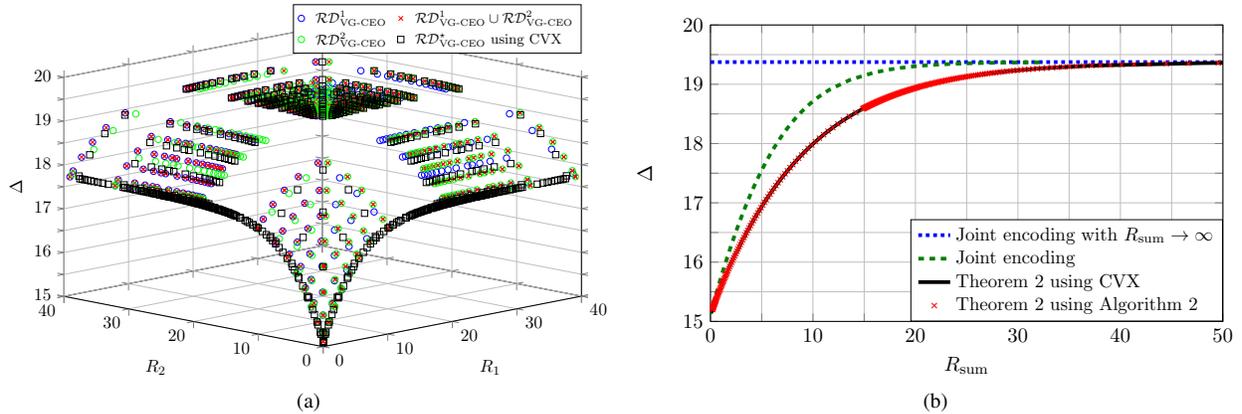


Fig. 5: Rate-information region of the vector Gaussian CEO network of Example 2. Numerical values are  $n_x = 3$  and  $n_0 = n_1 = n_2 = 4$ . (a): set of  $(R_1, R_2, \Delta)$  triples such  $(R_1, R_2, h(\mathbf{X}) - \Delta) \in \mathcal{RD}_{\text{VG-CEO}}^1 \cup \mathcal{RD}_{\text{VG-CEO}}^2$ , computed using Algorithm 2. (b): set of  $(R_{\text{sum}}, \Delta)$  pairs such  $R_{\text{sum}} = R_1 + R_2$  for some  $(R_1, R_2)$  for which  $(R_1, R_2, h(\mathbf{X}) - \Delta) \in \mathcal{RD}_{\text{VG-CEO}}^1 \cup \mathcal{RD}_{\text{VG-CEO}}^2$ .

**Example 2.** Consider an instance of the memoryless vector Gaussian CEO problem as described by (3) and (4) obtained by setting  $K = 2$ ,  $n_x = 3$  and  $n_0 = n_1 = n_2 = 4$ . We use Algorithm 2 to numerically approximate the set of  $(R_1, R_2, \Delta)$  triples such  $(R_1, R_2, h(\mathbf{X}) - \Delta)$  is in the union of the achievable regions  $\mathcal{RD}_{\text{VG-CEO}}^1$  and  $\mathcal{RD}_{\text{VG-CEO}}^2$ . The result is depicted in Figure 5a. The figure also shows the set of  $(R_1, R_2, \Delta)$  triples such  $(R_1, R_2, h(\mathbf{X}) - \Delta)$  lies in the region given by Theorem 2 evaluated for the example at hand. Figure 5b shows the

<sup>2</sup>We remind the reader that, as already mentioned, Algorithm 1 only converges to stationary points of the rate-distortion region.

set of  $(R_{\text{sum}}, \Delta)$  pairs such  $R_{\text{sum}} := R_1 + R_2$  for some  $(R_1, R_2)$  for which  $(R_1, R_2, h(\mathbf{X}) - \Delta)$  is in the union of  $\mathcal{RD}_{\text{VG-CEO}}^1$  and  $\mathcal{RD}_{\text{VG-CEO}}^2$ . The region is computed using two different approaches: i) using Algorithm 2 and ii) by directly evaluating the region obtained from Theorem 2 using the CVX optimization tool to find the maximizing covariances matrices  $(\mathbf{\Omega}_1, \mathbf{\Omega}_2)$  (note that this problem is convex and so CVX finds the optimal solution). It is worth-noting that Algorithm 2 converges to the optimal solution for the studied vector Gaussian CEO example, as is visible from the figure. For comparisons reasons, the figure also shows the performance of centralized or joint encoding, i.e., the case both agents observe both  $\mathbf{Y}_1$  and  $\mathbf{Y}_2$ ,

$$\Delta(R_{\text{sum}}) = \max_{P_{U|\mathbf{Y}_1, \mathbf{Y}_2} : I(U; \mathbf{Y}_1, \mathbf{Y}_2 | \mathbf{Y}_0) \leq R_{\text{sum}}} I(U, \mathbf{Y}_0; \mathbf{X}) . \quad (39)$$

Finally, we note that the information/sum-rate function (39) can be seen an extension of Chechik et al. Gaussian Information Bottleneck [26] to the case of side information  $\mathbf{Y}_0$  at the decoder. Figure 5b shows the loss in terms of information/sum-rate that is incurred by restricting the encoders to operate separately, i.e., distributed Information Bottleneck with side information at decoder. ■

## APPENDIX I

### PROOF OF THEOREM 1

For convenience, consider first the DM version of the CEO problem with decoder side information under logarithmic loss of Figure 1. It is assumed that for all  $\mathcal{S} \subseteq \mathcal{K} := \{1, \dots, K\}$ ,

$$\mathbf{Y}_{\mathcal{S}} \text{---} (\mathbf{X}, \mathbf{Y}_0) \text{---} \mathbf{Y}_{\mathcal{S}^c} , \quad (40)$$

forms a Markov chain in that order. The definitions for this model are similar to Definition 1 and Definition 2 and are omitted for brevity. The rate-distortion region of this problem can be obtained readily by applying [10, Theorem 10], which provides the rate-distortion region of the model without side information at decoder, to the modified setting in which the remote source is  $\tilde{\mathbf{X}} = (\mathbf{X}, \mathbf{Y}_0)$ , another agent (agent  $K + 1$ ) observes  $\mathbf{Y}_{K+1} = \mathbf{Y}_0$  and communicates at large rate  $R_{K+1} = \infty$  with the CEO, which wishes to estimates  $\tilde{\mathbf{X}}$  to within average logarithmic distortion  $D$  and has no own side information stream<sup>3</sup>. More specifically, it is given by the union of the set of all non-negative tuples  $(R_1, \dots, R_K, D)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,

$$D + \sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} I(\mathbf{Y}_k; U_k | \mathbf{X}, \mathbf{Y}_0, Q) + H(\mathbf{X} | U_{\mathcal{S}^c}, \mathbf{Y}_0, Q) , \quad (41)$$

for some joint measure of the form  $P_{\mathbf{Y}_0, \mathbf{Y}_{\mathcal{K}}, \mathbf{X}}(\mathbf{y}_0, \mathbf{y}_{\mathcal{K}}, \mathbf{x}) P_Q(q) \prod_{k=1}^K P_{U_k | \mathbf{Y}_k, Q}(u_k | \mathbf{y}_k, q)$ .

Also, let us define for this model the rate-information region  $\mathcal{RT}_{\text{CEO}}^*$  as the closure of all rate-information tuples  $(R_1, \dots, R_K, \Delta)$  for which there exist a blocklength  $n$ , encoding functions  $\{\phi_k^{(n)}\}_{k=1}^K$  and a decoding function

<sup>3</sup>Note that for the modified CEO setting the agents' observations are conditionally independent given the remote source  $\tilde{\mathbf{X}}$ .

$\psi^{(n)}$  such that

$$\begin{aligned} R_k &\geq \frac{1}{n} \log M_k^{(n)}, \quad \text{for } k = 1, \dots, K, \\ \Delta &\leq \frac{1}{n} I(\mathbf{X}^n; \psi^{(n)}(\phi_1^{(n)}(\mathbf{Y}_1^n), \dots, \phi_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n)). \end{aligned}$$

It is easy to see that a characterization of  $\mathcal{R}\mathcal{I}_{\text{CEO}}^*$  can be obtained by using (41) and substituting distortion levels  $D$  therein with  $(\Delta := H(\mathbf{X}) - D)$ .

**Proposition 5.** *The rate-information region  $\mathcal{R}\mathcal{I}_{\text{CEO}}^*$  of the vector DM CEO problem under logarithmic loss is given by the set of all non-negative tuples  $(R_1, \dots, R_K, \Delta)$  that satisfy, for all subsets  $\mathcal{S} \subseteq \mathcal{K}$ ,*

$$\sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} I(\mathbf{Y}_k; U_k | \mathbf{X}, \mathbf{Y}_0, Q) - I(\mathbf{X}; U_{\mathcal{S}^c}, \mathbf{Y}_0, Q) + \Delta,$$

for some joint measure of the form  $P_{\mathbf{Y}_0, \mathbf{Y}_{\mathcal{K}}, \mathbf{X}}(\mathbf{y}_0, \mathbf{y}_{\mathcal{K}}, \mathbf{x}) P_Q(q) \prod_{k=1}^K P_{U_k | \mathbf{Y}_k, Q}(u_k | \mathbf{y}_k, q)$ . ■

The region  $\mathcal{R}\mathcal{I}_{\text{CEO}}^*$  involves mutual information terms only (not entropies); and, so, using a standard discretization argument, it can be easily shown that a characterization of this region in the case of continuous alphabets is also given by Proposition 5. It is well-known that the rate region of the DM CEO problem under logarithmic loss can also be used in the case of continuous alphabets (e.g., Gaussian sources [61], [62]).

Let us now return to the vector Gaussian CEO problem under logarithmic loss that we study in this section. First, we state the following lemma, whose proof is easy and is omitted for brevity.

**Lemma 5.**  $(R_1, \dots, R_K, D) \in \mathcal{R}\mathcal{D}_{\text{VG-CEO}}^*$  if and only if  $(R_1, \dots, R_K, h(\mathbf{X}) - D) \in \mathcal{R}\mathcal{I}_{\text{CEO}}^*$ . ■

Summarizing, using Proposition 5 and Lemma 5 it follows that  $\mathcal{R}\mathcal{D}_{\text{VG-CEO}}^* = \mathcal{R}\mathcal{D}_{\text{CEO}}^{\text{I}}$ . To complete the proof of Theorem 1, it remains to show that  $\mathcal{R}\mathcal{D}_{\text{CEO}}^{\text{I}} = \mathcal{R}\mathcal{D}_{\text{CEO}}^{\text{II}}$ ; and this follows by reasoning along the submodularity arguments of the proof of [10, Theorem 10].

## APPENDIX II

### PROOF OF CONVERSE OF THEOREM 2

The proof of Theorem 2 relies on deriving an outer bound on the region  $\mathcal{R}\mathcal{D}_{\text{CEO}}^{\text{I}}$  given by Theorem 1. In doing so, we use the technique of [18, Theorem 8] which relies on the de Bruijn identity and the properties of Fisher information; and extend the argument to account for the time-sharing variable  $Q$  and side information  $\mathbf{Y}_0$ .

We first state the following lemma.

**Lemma 6.** [18], [63] *Let  $(\mathbf{X}, \mathbf{Y})$  be a pair of random vectors with pmf  $p(\mathbf{x}, \mathbf{y})$ . We have*

$$\log |(\pi e) \mathbf{J}^{-1}(\mathbf{X} | \mathbf{Y})| \leq h(\mathbf{X} | \mathbf{Y}) \leq \log |(\pi e) \text{mmse}(\mathbf{X} | \mathbf{Y})|,$$

where the conditional Fisher information matrix is defined as

$$\mathbf{J}(\mathbf{X}|\mathbf{Y}) := \mathbb{E}[\nabla \log p(\mathbf{X}|\mathbf{Y}) \nabla \log p(\mathbf{X}|\mathbf{Y})^\dagger],$$

and the minimum mean squared error (MMSE) matrix is

$$\text{mmse}(\mathbf{X}|\mathbf{Y}) := \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}])(\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}])^\dagger]. \quad \blacksquare$$

Now, we derive an outer bound on (7) as follows. For each  $q \in \mathcal{Q}$  and fixed pmf  $\prod_{k=1}^K p(u_k|\mathbf{y}_k, q)$ , choose  $\{\boldsymbol{\Omega}_{k,q}\}_{k=1}^K$  satisfying  $\mathbf{0} \preceq \boldsymbol{\Omega}_{k,q} \preceq \boldsymbol{\Sigma}_k^{-1}$  such that

$$\text{mmse}(\mathbf{Y}_k|\mathbf{X}, U_{k,q}, \mathbf{Y}_0, q) = \boldsymbol{\Sigma}_k - \boldsymbol{\Sigma}_k \boldsymbol{\Omega}_{k,q} \boldsymbol{\Sigma}_k. \quad (42)$$

Such  $\boldsymbol{\Omega}_{k,q}$  always exists since, for all  $q \in \mathcal{Q}$ ,  $k \in \mathcal{K}$ , we have

$$\mathbf{0} \preceq \text{mmse}(\mathbf{Y}_k|\mathbf{X}, U_{k,q}, \mathbf{Y}_0, q) \preceq \boldsymbol{\Sigma}_{\mathbf{y}_k|\mathbf{x}, \mathbf{y}_0} = \boldsymbol{\Sigma}_{\mathbf{n}_k|\mathbf{n}_0} = \boldsymbol{\Sigma}_k.$$

Then, for  $k \in \mathcal{K}$  and  $q \in \mathcal{Q}$ , we have

$$\begin{aligned} I(\mathbf{Y}_k; U_k|\mathbf{X}, \mathbf{Y}_0, Q = q) &= \log |(\pi e) \boldsymbol{\Sigma}_k| - h(\mathbf{Y}_k|\mathbf{X}, U_{k,q}, \mathbf{Y}_0, Q = q) \\ &\stackrel{(a)}{\geq} \log |\boldsymbol{\Sigma}_k| - \log |\text{mmse}(\mathbf{Y}_k|\mathbf{X}, U_{k,q}, \mathbf{Y}_0, Q = q)| \\ &\stackrel{(b)}{=} -\log |\mathbf{I} - \boldsymbol{\Omega}_{k,q} \boldsymbol{\Sigma}_k|, \end{aligned} \quad (43)$$

where (a) is due to Lemma 6; and (b) is due to (42).

For convenience, the matrix  $\boldsymbol{\Lambda}_{\bar{\mathcal{S}},q}$  is defined as follows

$$\boldsymbol{\Lambda}_{\mathcal{S},q} := \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{diag}(\{\boldsymbol{\Sigma}_k - \boldsymbol{\Sigma}_k \boldsymbol{\Omega}_{k,q} \boldsymbol{\Sigma}_k\}_{k \in \mathcal{S}^c}) \end{bmatrix}. \quad (44)$$

Then, for  $q \in \mathcal{Q}$  and  $\mathcal{S} \subseteq \mathcal{K}$ , we have

$$\begin{aligned} h(\mathbf{X}|U_{\mathcal{S}^c,q}, \mathbf{Y}_0, Q = q) &\stackrel{(a)}{\geq} \log |(\pi e) \mathbf{J}^{-1}(\mathbf{X}|U_{\mathcal{S}^c,q}, \mathbf{Y}_0, q)| \\ &\stackrel{(b)}{=} \log \left| (\pi e) \left( \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{H}_{\bar{\mathcal{S}}}^\dagger \boldsymbol{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1} (\mathbf{I} - \boldsymbol{\Lambda}_{\bar{\mathcal{S}},q} \boldsymbol{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1}) \mathbf{H}_{\bar{\mathcal{S}}} \right)^{-1} \right|, \end{aligned} \quad (45)$$

where (a) follows from Lemma 6; and for (b), we use the connection of the MMSE and the Fisher information to show the following equality

$$\mathbf{J}(\mathbf{X}|U_{\mathcal{S}^c,q}, \mathbf{Y}_0, q) = \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} + \mathbf{H}_{\bar{\mathcal{S}}}^\dagger \boldsymbol{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1} (\mathbf{I} - \boldsymbol{\Lambda}_{\bar{\mathcal{S}},q} \boldsymbol{\Sigma}_{\mathbf{n}_{\bar{\mathcal{S}}}}^{-1}) \mathbf{H}_{\bar{\mathcal{S}}}. \quad (46)$$

In order to proof (46), we use de Bruijn identity to relate the Fisher information with the MMSE as given in the following lemma.

**Lemma 7.** [18], [64] Let  $(\mathbf{V}_1, \mathbf{V}_2)$  be a random vector with finite second moments and  $\mathbf{Z} \sim \mathcal{CN}(\mathbf{0}, \Sigma_{\mathbf{z}})$  independent of  $(\mathbf{V}_1, \mathbf{V}_2)$ . Then

$$\text{mmse}(\mathbf{V}_2 | \mathbf{V}_1, \mathbf{V}_2 + \mathbf{Z}) = \Sigma_{\mathbf{z}} - \Sigma_{\mathbf{z}} \mathbf{J}(\mathbf{V}_2 + \mathbf{Z} | \mathbf{V}_1) \Sigma_{\mathbf{z}}. \quad \blacksquare$$

From MMSE estimation of Gaussian random vectors, for  $\mathcal{S} \subseteq \mathcal{K}$ , we have

$$\mathbf{X} = \mathbb{E}[\mathbf{X} | \mathbf{Y}_{\mathcal{S}}] + \mathbf{W}_{\mathcal{S}} = \mathbf{G}_{\mathcal{S}} \mathbf{Y}_{\mathcal{S}} + \mathbf{W}_{\mathcal{S}}, \quad (47)$$

where  $\mathbf{G}_{\mathcal{S}} := \Sigma_{\mathbf{w}_{\mathcal{S}}} \mathbf{H}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1}$ , and  $\mathbf{W}_{\mathcal{S}} \sim \mathcal{CN}(\mathbf{0}, \Sigma_{\mathbf{w}_{\mathcal{S}}})$  is a Gaussian vector that is independent of  $\mathbf{Y}_{\mathcal{S}}$  and

$$\Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} := \Sigma_{\mathbf{x}}^{-1} + \mathbf{H}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1} \mathbf{H}_{\mathcal{S}}. \quad (48)$$

Now we show that the cross-terms of  $\text{mmse}(\mathbf{Y}_{\mathcal{S}^c} | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q)$  are zero (similarly to [18, Appendix V]). For  $i \in \mathcal{S}^c$  and  $j \neq i$ , we have

$$\begin{aligned} & \mathbb{E}[(Y_i - \mathbb{E}[Y_i | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q])(Y_j - \mathbb{E}[Y_j | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q])^{\dagger}] \\ & \stackrel{(a)}{=} \mathbb{E} \left[ \mathbb{E}[(Y_i - \mathbb{E}[Y_i | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q])(Y_j - \mathbb{E}[Y_j | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q])^{\dagger} | \mathbf{X}, \mathbf{Y}_0] \right] \\ & \stackrel{(b)}{=} \mathbb{E} \left[ \mathbb{E}[(Y_i - \mathbb{E}[Y_i | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q]) | \mathbf{X}, \mathbf{Y}_0] \mathbb{E}[(Y_j - \mathbb{E}[Y_j | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q])^{\dagger} | \mathbf{X}, \mathbf{Y}_0] \right] = \mathbf{0}, \end{aligned} \quad (49)$$

where (a) is due to the law of total expectation; (b) is due to the Markov chain  $\mathbf{Y}_k \text{---} (\mathbf{X}, \mathbf{Y}_0) \text{---} \mathbf{Y}_{\mathcal{K} \setminus k}$ .

Then, for  $k \in \mathcal{K}$  and  $q \in \mathcal{Q}$ , we have

$$\begin{aligned} \text{mmse}(\mathbf{G}_{\mathcal{S}} \mathbf{Y}_{\mathcal{S}} | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q) &= \mathbf{G}_{\mathcal{S}} \text{mmse}(\mathbf{Y}_{\mathcal{S}} | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q) \mathbf{G}_{\mathcal{S}}^{\dagger} \\ & \stackrel{(a)}{=} \mathbf{G}_{\mathcal{S}} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{diag}(\{\text{mmse}(\mathbf{Y}_k | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q)\}_{k \in \mathcal{S}^c}) \end{bmatrix} \mathbf{G}_{\mathcal{S}}^{\dagger} \\ & \stackrel{(b)}{=} \mathbf{G}_{\mathcal{S}} \Lambda_{\mathcal{S}, q} \mathbf{G}_{\mathcal{S}}^{\dagger}, \end{aligned} \quad (50)$$

where (a) follows since the cross-terms are zero as shown in (49); and (b) follows due to (42) and the definition of  $\Lambda_{\mathcal{S}, q}$  given in (44).

Finally, we obtain the equality (46) by applying Lemma 7 and noting (47) as follows

$$\begin{aligned} \mathbf{J}(\mathbf{X} | U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q) & \stackrel{(a)}{=} \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} - \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} \text{mmse}(\mathbf{G}_{\mathcal{S}} \mathbf{Y}_{\mathcal{S}} | \mathbf{X}, U_{\mathcal{S}^c, q}, \mathbf{Y}_0, q) \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} \\ & \stackrel{(b)}{=} \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} - \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} \mathbf{G}_{\mathcal{S}} \Lambda_{\mathcal{S}, q} \mathbf{G}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1} \\ & \stackrel{(c)}{=} \Sigma_{\mathbf{x}}^{-1} + \mathbf{H}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1} \mathbf{H}_{\mathcal{S}} - \mathbf{H}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1} \Lambda_{\mathcal{S}, q} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1} \mathbf{H}_{\mathcal{S}} \\ & = \Sigma_{\mathbf{x}}^{-1} + \mathbf{H}_{\mathcal{S}}^{\dagger} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1} (\mathbf{I} - \Lambda_{\mathcal{S}, q} \Sigma_{\mathbf{n}_{\mathcal{S}}}^{-1}) \mathbf{H}_{\mathcal{S}}, \end{aligned}$$

where (a) is due to Lemma 7; (b) is due to (50); and (c) follows due to the definitions of  $\Sigma_{\mathbf{w}_{\mathcal{S}}}^{-1}$  and  $\mathbf{G}_{\mathcal{S}}$ .

Next, we average (43) and (45) over the time-sharing  $Q$  and letting  $\mathbf{\Omega}_k := \sum_{q \in \mathcal{Q}} p(q) \mathbf{\Omega}_{k,q}$ , we obtain the lower bound

$$\begin{aligned}
I(\mathbf{Y}_k; \mathbf{U}_k | \mathbf{X}, \mathbf{Y}_0, Q) &= \sum_{q \in \mathcal{Q}} p(q) I(\mathbf{Y}_k; \mathbf{U}_k | \mathbf{X}, \mathbf{Y}_0, Q = q) \\
&\stackrel{(a)}{\geq} - \sum_{q \in \mathcal{Q}} p(q) \log |\mathbf{I} - \mathbf{\Omega}_{k,q} \mathbf{\Sigma}_k| \\
&\stackrel{(b)}{\geq} - \log |\mathbf{I} - \sum_{q \in \mathcal{Q}} p(q) \mathbf{\Omega}_{k,q} \mathbf{\Sigma}_k| \\
&= - \log |\mathbf{I} - \mathbf{\Omega}_k \mathbf{\Sigma}_k|, \tag{51}
\end{aligned}$$

where (a) follows from (43); and (b) follows from the concavity of the log-det function and Jensen's Inequality. Besides, we can derive the following lower bound

$$\begin{aligned}
h(\mathbf{X} | U_{S^c}, \mathbf{Y}_0, Q) &= \sum_{q \in \mathcal{Q}} p(q) h(\mathbf{X} | U_{S^c, q}, \mathbf{Y}_0, Q = q) \\
&\stackrel{(a)}{\geq} \sum_{q \in \mathcal{Q}} p(q) \log \left| (\pi e) \left( \mathbf{\Sigma}_x^{-1} + \mathbf{H}_S^\dagger \mathbf{\Sigma}_{n_S}^{-1} (\mathbf{I} - \mathbf{\Lambda}_{S,q} \mathbf{\Sigma}_{n_S}^{-1}) \mathbf{H}_S \right)^{-1} \right| \\
&\stackrel{(b)}{\geq} \log \left| (\pi e) \left( \mathbf{\Sigma}_x^{-1} + \mathbf{H}_S^\dagger \mathbf{\Sigma}_{n_S}^{-1} (\mathbf{I} - \mathbf{\Lambda}_S \mathbf{\Sigma}_{n_S}^{-1}) \mathbf{H}_S \right)^{-1} \right|, \tag{52}
\end{aligned}$$

where (a) is due to (45); and (b) is due to the concavity of the log-det function and Jensen's inequality and the definition of  $\mathbf{\Lambda}_S$  given in (11).

Finally, the outer bound on  $\mathcal{RD}_{\text{VG-CEO}}^*$  is obtained by applying (51) and (52) in (7), noting that  $\mathbf{\Omega}_k = \sum_{q \in \mathcal{Q}} p(q) \mathbf{\Omega}_{k,q} \preceq \mathbf{\Sigma}_k^{-1}$  since  $\mathbf{0} \preceq \mathbf{\Omega}_{k,q} \preceq \mathbf{\Sigma}_k^{-1}$ , and taking the union over  $\mathbf{\Omega}_k$  satisfying  $\mathbf{0} \preceq \mathbf{\Omega}_k \preceq \mathbf{\Sigma}_k^{-1}$ .

### APPENDIX III

#### PROOF OF THEOREM 3

We first present the following lemma, which essentially states that Theorem 2 provides an outer bound on  $\mathcal{RD}_{\text{VG-CEO}}^{\text{det}}$ .

**Lemma 8.** *If  $(R_1, \dots, R_K, D) \in \mathcal{RD}_{\text{VG-CEO}}^{\text{det}}$ , then  $(R_1, \dots, R_K, \log(\pi e)^{n_x} D) \in \mathcal{RD}_{\text{CEO}}^{\text{I}}$ .*

*Proof.* Let a tuple  $(R_1, \dots, R_K, D) \in \mathcal{RD}_{\text{VG-CEO}}^{\text{det}}$  be given. Then, there exist a blocklength  $n$ ,  $K$  encoding functions  $\{\check{\phi}_k^{(n)}\}_{k=1}^K$  and a decoding function  $\check{\psi}^{(n)}$  such that

$$\begin{aligned}
R_k &\geq \frac{1}{n} \log M_k^{(n)}, \quad \text{for } k = 1, \dots, K, \\
D &\geq \left| \frac{1}{n} \sum_{i=1}^n \text{mmse}(\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \right|. \tag{53}
\end{aligned}$$

We need to show that there exist  $(U_1, \dots, U_K, Q)$  such that

$$\sum_{k \in \mathcal{S}} R_k + \log(\pi e)^{n_x} D \geq \sum_{k \in \mathcal{S}} I(\mathbf{Y}_k; U_k | \mathbf{X}, \mathbf{Y}_0, Q) + h(\mathbf{X} | U_{\mathcal{S}^c}, \mathbf{Y}_0, Q), \quad \text{for } \mathcal{S} \subseteq \mathcal{K}. \quad (54)$$

Let us define

$$\bar{\Delta}^{(n)} := \frac{1}{n} h(\mathbf{X}^n | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n).$$

It is easy to justify that expected distortion  $\bar{\Delta}^{(n)}$  is achievable under logarithmic loss (see Theorem 1). Then, following straightforwardly the lines in the proof of [10, Theorem 10], we have

$$\sum_{k \in \mathcal{S}} R_k \geq \sum_{k \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n I(\mathbf{Y}_{k,i}; U_{k,i} | \mathbf{X}_i, \mathbf{Y}_{0,i}, Q_i) + \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i | U_{\mathcal{S}^c,i}, \mathbf{Y}_{0,i}, Q_i) - \bar{\Delta}^{(n)}. \quad (55)$$

Next, we upper bound  $\bar{\Delta}^{(n)}$  in terms of  $D$  as follows

$$\begin{aligned} \bar{\Delta}^{(n)} &= \frac{1}{n} h(\mathbf{X}^n | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \\ &= \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i | \mathbf{X}_{i+1}^n, \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \\ &= \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i | J_{\mathcal{K}}] | \mathbf{X}_{i+1}^n, \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \\ &\stackrel{(a)}{\leq} \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n]) \\ &\stackrel{(b)}{\leq} \frac{1}{n} \sum_{i=1}^n \log(\pi e)^{n_x} \left| \text{mmse}(\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \right| \\ &\stackrel{(c)}{\leq} \log(\pi e)^{n_x} \left| \frac{1}{n} \sum_{i=1}^n \text{mmse}(\mathbf{X}_i | \check{\phi}_1^{(n)}(\mathbf{Y}_1^n), \dots, \check{\phi}_K^{(n)}(\mathbf{Y}_K^n), \mathbf{Y}_0^n) \right| \\ &\stackrel{(d)}{\leq} \log(\pi e)^{n_x} D, \end{aligned} \quad (56)$$

where (a) holds since conditioning reduces entropy; (b) is due to the maximal differential entropy lemma; (c) is due to the convexity of the log-det function and Jensen's inequality; and (d) is due to (53).

Combining (56) with (55), and using standard arguments for single-letterization, we get (54); and this completes the proof of the lemma.  $\square$

The proof of Theorem 3 is as follows. By Lemma 8 and Proposition 1, there must exist Gaussian test channels  $(V_1^G, \dots, V_K^G)$  and a time-sharing random variable  $Q'$ , with joint distribution that factorizes as

$$P_{\mathbf{X}, \mathbf{Y}_0}(\mathbf{x}, \mathbf{y}_0) \prod_{k=1}^K P_{\mathbf{Y}_k | \mathbf{X}, \mathbf{Y}_0}(\mathbf{y}_k | \mathbf{x}, \mathbf{y}_0) P_{Q'}(q') \prod_{k=1}^K P_{V_k | \mathbf{Y}_k, Q'}(v_k | \mathbf{y}_k, q'),$$

such that the following holds

$$\sum_{k \in \mathcal{S}} R_k \geq I(\mathbf{Y}_{\mathcal{S}}; V_{\mathcal{S}}^G | V_{\mathcal{S}^c}^G, \mathbf{Y}_0, Q'), \quad \text{for } \mathcal{S} \subseteq \mathcal{K}, \quad (57)$$

$$\log(\pi e)^{n_x} D \geq h(\mathbf{X} | V_1^G, \dots, V_K^G, \mathbf{Y}_0, Q'). \quad (58)$$

This is clearly achievable by the Berger-Tung coding scheme with Gaussian test channels and time-sharing  $Q'$ , since the achievable error matrix under quadratic distortion has determinant that satisfies

$$\log((\pi e)^{n_x} |\text{mmse}(\mathbf{X} | V_1^G, \dots, V_K^G, \mathbf{Y}_0, Q')|) = h(\mathbf{X} | V_1^G, \dots, V_K^G, \mathbf{Y}_0, Q').$$

The above shows that the rate-distortion region of the quadratic vector Gaussian CEO problem with determinant constraint is given by (58), i.e.,  $\mathcal{RD}_{\text{CEO}}^{\text{II}}$  (with distortion parameter  $\log(\pi e)^{n_x} D$ ). Recalling that  $\mathcal{RD}_{\text{CEO}}^{\text{II}} = \mathcal{RD}_{\text{CEO}}^{\text{I}} = \mathcal{RD}_{\text{VG-CEO}}^*$ , and substituting in Theorem 2 using distortion level  $\log(\pi e)^{n_x} D$  completes the proof.

#### APPENDIX IV

##### PROOF OF PROPOSITION 2

Let a triple mappings  $(Q_{U_1|Y_1}, Q_{U_2|Y_2}, Q_{\hat{X}|U_1, U_2, Y_0})$  be given. It is easy to see that the probability of classification error of the classifier  $Q_{\hat{X}|Y_0, Y_1, Y_2}$  as defined by (20) satisfies

$$P_{\mathcal{E}}(Q_{\hat{X}|Y_0, Y_1, Y_2}) \leq \mathbb{E}_{P_{X, Y_0, Y_1, Y_2}}[-\log Q_{\hat{X}|Y_0, Y_1, Y_2}(X|Y_0, Y_1, Y_2)]. \quad (59)$$

Applying Jensen's inequality on the right hand side (RHS) of (59), using the concavity of the logarithm function, and combining with the fact that the exponential function increases monotonically, the probability of classification error can be further bounded as

$$P_{\mathcal{E}}(Q_{\hat{X}|Y_0, Y_1, Y_2}) \leq 1 - \exp\left(-\mathbb{E}_{P_{X, Y_0, Y_1, Y_2}}[-\log Q_{\hat{X}|Y_0, Y_1, Y_2}(X|Y_0, Y_1, Y_2)]\right). \quad (60)$$

Using (19) and continuing from (60), we get

$$\begin{aligned} P_{\mathcal{E}}(Q_{\hat{X}|Y_0, Y_1, Y_2}) &\leq 1 - \exp\left(-\mathbb{E}_{P_{X, Y_0, Y_1, Y_2}}[-\log \mathbb{E}_{Q_{U_1|Y_1}} \mathbb{E}_{Q_{U_2|Y_2}} [Q_{\hat{X}|U_1, U_2, Y_0}(X|U_1, U_2, Y_0)]]\right) \\ &\leq 1 - \exp\left(-\mathbb{E}_{P_{X, Y_0, Y_1, Y_2}} \mathbb{E}_{Q_{U_1|Y_1}} \mathbb{E}_{Q_{U_2|Y_2}} [-\log [Q_{\hat{X}|U_1, U_2, Y_0}(X|U_1, U_2, Y_0)]]\right), \end{aligned} \quad (61)$$

where the last inequality follows by applying Jensen's inequality and using the concavity of the logarithm function. Noticing that the term in the exponential function in the RHS of (61),

$$\mathcal{D}(Q_{U_1|Y_1}, Q_{U_2|Y_2}, Q_{\hat{X}|U_1, U_2, Y_0}) := \mathbb{E}_{P_{X, Y_0, Y_1, Y_2}} \mathbb{E}_{Q_{U_1|Y_1}} \mathbb{E}_{Q_{U_2|Y_2}} [-\log Q_{\hat{X}|U_1, U_2, Y_0}(X|U_1, U_2, Y_0)], \quad (62)$$

is the average logarithmic loss, or cross-entropy risk, of the triple  $(Q_{U_1|Y_1}, Q_{U_2|Y_2}, Q_{\hat{X}|U_1, U_2, Y_0})$ ; the inequality (61) implies that minimizing the average logarithmic loss distortion leads to classifier with smaller (bound on) its classification error. Using Theorem 1, the minimum average logarithmic loss, minimized over all mappings

$Q_{U_1|Y_1} : \mathcal{Y}_1 \rightarrow \mathcal{P}(\mathcal{U}_1)$  and  $Q_{U_2|Y_2} : \mathcal{Y}_2 \rightarrow \mathcal{P}(\mathcal{U}_2)$  that have description lengths no more than  $R_1$  and  $R_2$  bits per-sample, respectively, as well as all choices of  $Q_{\hat{X}|U_1, U_2, Y_0} : \mathcal{U}_1 \times \mathcal{U}_2 \times \mathcal{Y}_0 \rightarrow \mathcal{P}(\mathcal{X})$ , is

$$D^*(R_1, R_2) = \inf\{D : (R_1, R_2, D) \in \mathcal{RD}_{\text{CEO}}^*\}. \quad (63)$$

Thus, the direct part of Theorem 1 guarantees the existence of a classifier  $Q_{\hat{X}|Y_0, Y_1, Y_2}^*$  whose probability of error satisfies the bound given in Proposition 2.

#### APPENDIX V PROOF OF LEMMA 2

First, we rewrite  $F_s(\mathbf{P})$  in (28). To that end, the second term of the RHS of (28) can be proceeded as

$$\begin{aligned} I(Y_1; U_1|U_2, Y_0) &\stackrel{(a)}{=} I(X, Y_1; U_1|U_2, Y_0) = I(X; U_1|U_2, Y_0) + I(Y_1; U_1|U_2, Y_0, X) \\ &\stackrel{(b)}{=} I(X; U_1|U_2, Y_0) + I(Y_1; U_1|X, Y_0) \\ &= I(X; U_1|U_2, Y_0) + I(Y_1, X; U_1|Y_0) - I(X; U_1|Y_0) \\ &\stackrel{(c)}{=} I(X; U_1|U_2, Y_0) + I(Y_1; U_1|Y_0) - I(X; U_1|Y_0) \\ &= H(X|U_2, Y_0) - H(X|U_1, U_2, Y_0) + H(U_1|Y_0) - H(U_1|Y_0, Y_1) - H(X|Y_0) + H(X|U_1, Y_0) \\ &= H(X|U_2, Y_0) - H(X|U_1, U_2, Y_0) + H(U_1) - H(Y_0) + H(Y_0|U_1) \\ &\quad - H(U_1|Y_0, Y_1) - H(X|Y_0) + H(X|U_1, Y_0), \end{aligned} \quad (64)$$

and, the third term of the RHS of (28) can be written as

$$\begin{aligned} I(Y_2; U_2|Y_0) &= H(U_2|Y_0) - H(U_2|Y_0, Y_2) \stackrel{(d)}{=} H(U_2|Y_0) - H(U_2|Y_2) \\ &= H(U_2) - H(Y_0) + H(Y_0|U_2) - H(U_2|Y_2), \end{aligned} \quad (65)$$

where (a), (b), (c) and (d) follows due to the Markov chain  $U_1 \text{---} Y_1 \text{---} (X, Y_0) \text{---} Y_2 \text{---} U_2$ .

By applying (64) and (65) in (28), we have

$$\begin{aligned} F_s(\mathbf{P}) &= -s_1 H(X|Y_0) - (s_1 + s_2)H(Y_0) + (1 - s_1)H(X|U_1, U_2, Y_0) \\ &\quad + s_1 H(X|U_1, Y_0) + s_1 H(X|U_2, Y_0) + s_1 H(U_1) - s_1 H(U_1|Y_1) \\ &\quad + s_2 H(U_2) - s_2 H(U_2|Y_2) + s_1 H(Y_0|U_1) + s_2 H(Y_0|U_2) \\ &= -s_1 H(X|Y_0) - (s_1 + s_2)H(Y_0) \\ &\quad - (1 - s_1) \sum_{u_1 u_2 x y_0} p(u_1, u_2, x, y_0) \log p(x|u_1, u_2, y_0) \\ &\quad - s_1 \sum_{u_1 x y_0} p(u_1, x, y_0) \log p(x|u_1, y_0) - s_1 \sum_{u_2 x y_0} p(u_2, x, y_0) \log p(x|u_2, y_0) \end{aligned}$$

$$\begin{aligned}
& -s_1 \sum_{u_1} p(u_1) \log p(u_1) + s_1 \sum_{u_1 y_1} p(u_1, y_1) \log p(u_1|y_1) \\
& -s_2 \sum_{u_2} p(u_2) \log p(u_2) + s_2 \sum_{u_2 y_2} p(u_2, y_2) \log p(u_2|y_2) \\
& -s_1 \sum_{u_1 y_0} p(u_1, y_0) \log p(y_0|u_1) - s_2 \sum_{u_2 y_0} p(u_2, y_0) \log p(y_0|u_2), \tag{66}
\end{aligned}$$

Then, marginalizing (66) over variables  $X, Y_0, Y_1, Y_2$ , and using the Markov chain  $U_1 \ominus Y_1 \ominus (X, Y_0) \ominus Y_2 \ominus U_2$ , it is easy to see that  $F_{\mathbf{s}}(\mathbf{P})$  can be written as

$$\begin{aligned}
F_{\mathbf{s}}(\mathbf{P}) &= -s_1 H(X|Y_0) - (s_1 + s_2) H(Y_0) \\
&+ \mathbb{E}_{P_{X, Y_0, Y_1, Y_2}} \left[ (1 - s_1) \mathbb{E}_{P_{U_1|Y_1}} \mathbb{E}_{P_{U_2|Y_2}} [-\log P_{X|U_1, U_2, Y_0}] \right. \\
&\quad + s_1 \mathbb{E}_{P_{U_1|Y_1}} [-\log P_{X|U_1, Y_0}] + s_1 \mathbb{E}_{P_{U_2|Y_2}} [-\log P_{X|U_2, Y_0}] \\
&\quad + s_1 D_{\text{KL}}(P_{U_1|Y_1} \| P_{U_1}) + s_2 D_{\text{KL}}(P_{U_2|Y_2} \| P_{U_2}) \\
&\quad \left. + s_1 \mathbb{E}_{P_{U_1|Y_1}} [-\log P_{Y_0|U_1}] + s_2 \mathbb{E}_{P_{U_2|Y_2}} [-\log P_{Y_0|U_2}] \right]. \tag{67}
\end{aligned}$$

Hence, we have

$$\begin{aligned}
F_{\mathbf{s}}(\mathbf{P}, \mathbf{Q}) - F_{\mathbf{s}}(\mathbf{P}) &= (1 - s_1) \mathbb{E}_{U_1, U_2, Y_0} [D_{\text{KL}}(P_{X|U_1, U_2, Y_0} \| Q_{X|U_1, U_2, Y_0})] \\
&+ s_1 \mathbb{E}_{U_1, Y_0} [D_{\text{KL}}(P_{X|U_1, Y_0} \| Q_{X|U_1, Y_0})] + s_1 \mathbb{E}_{U_2, Y_0} [D_{\text{KL}}(P_{X|U_2, Y_0} \| Q_{X|U_2, Y_0})] \\
&+ s_1 D_{\text{KL}}(P_{U_1} \| Q_{U_1}) + s_2 D_{\text{KL}}(P_{U_2} \| Q_{U_2}) \\
&+ s_1 \mathbb{E}_{U_1} [D_{\text{KL}}(P_{Y_0|U_1} \| Q_{Y_0|U_1})] + s_2 \mathbb{E}_{U_2} [D_{\text{KL}}(P_{Y_0|U_2} \| Q_{Y_0|U_2})] \geq 0,
\end{aligned}$$

where it holds with equality if and only if (32) is satisfied. Note that we have the relation  $1 - s_1 \geq 0$  due to Lemma 1. This completes the proof.

## APPENDIX VI

### PROOF OF LEMMA 4

We have that  $F_{\mathbf{s}}(\mathbf{P}, \mathbf{Q})$  is convex in  $\mathbf{P}$  from Lemma 3. For a given  $\mathbf{Q}$  and  $\mathbf{s}$ , in order to minimize  $F_{\mathbf{s}}(\mathbf{P}, \mathbf{Q})$  over the convex set of pmfs  $\mathbf{P}$ , let us define the Lagrangian as

$$\mathcal{L}_{\mathbf{s}}(\mathbf{P}, \mathbf{Q}, \boldsymbol{\lambda}) := F_{\mathbf{s}}(\mathbf{P}, \mathbf{Q}) + \sum_{y_1} \lambda_1(y_1) [1 - \sum_{u_1} p(u_1|y_1)] + \sum_{y_2} \lambda_2(y_2) [1 - \sum_{u_2} p(u_2|y_2)],$$

where  $\lambda_1(y_1) \geq 0$  and  $\lambda_2(y_2) \geq 0$  are the Lagrange multipliers corresponding the constrains  $\sum_{u_k} p(u_k|y_k) = 1$ ,  $y_k \in \mathcal{Y}_k$ ,  $k = 1, 2$ , of the pmfs  $P_{U_1|Y_1}$  and  $P_{U_2|Y_2}$ , respectively. Due to the convexity of  $F_{\mathbf{s}}(\mathbf{P}, \mathbf{Q})$ , the KKT conditions are necessary and sufficient for optimality. By applying the KKT conditions

$$\frac{\partial \mathcal{L}_{\mathbf{s}}(\mathbf{P}, \mathbf{Q}, \boldsymbol{\lambda})}{\partial p(u_1|y_1)} = 0, \quad \frac{\partial \mathcal{L}_{\mathbf{s}}(\mathbf{P}, \mathbf{Q}, \boldsymbol{\lambda})}{\partial p(u_2|y_2)} = 0,$$

and arranging terms, we obtain

$$\begin{aligned}
& \log p(u_k|y_k) \\
&= \log q(u_k) + \frac{1-s_1}{s_k} \sum_{u_{\bar{k}} x y_0} p(x, y_0|y_k) p(u_{\bar{k}}|x, y_0) \log q(x|u_k, u_{\bar{k}}, y_0) \\
&\quad + \frac{s_1}{s_k} \sum_{x y_0} p(x, y_0|y_k) \log q(x|u_k, y_0) + \sum_{y_0} p(y_0|y_k) \log q(y_0|u_k) + \frac{\lambda_k(y_k)}{s_k p(y_k)} - 1 \\
&= \log q(u_k) + \frac{1-s_1}{s_k} \sum_{u_{\bar{k}} y_0} p(u_{\bar{k}}, y_0|y_k) \sum_x p(x|y_k, u_{\bar{k}}, y_0) \log q(x|u_k, u_{\bar{k}}, y_0) \\
&\quad + \frac{s_1}{s_k} \sum_{y_0} p(y_0|y_k) \sum_x p(x|y_k, y_0) \log q(x|u_k, y_0) + \sum_{y_0} p(y_0|y_k) \log q(y_0|u_k) + \frac{\lambda_k(y_k)}{s_k p(y_k)} - 1 \\
&= \log q(u_k) - \frac{1-s_1}{s_k} \sum_{u_{\bar{k}} y_0} p(u_{\bar{k}}, y_0|y_k) \sum_x p(x|y_k, u_{\bar{k}}, y_0) \log \frac{p(x|y_k, u_{\bar{k}}, y_0)}{q(x|u_k, u_{\bar{k}}, y_0)} \frac{1}{p(x|y_k, u_{\bar{k}}, y_0)} + \frac{\lambda_k(y_k)}{s_k p(y_k)} - 1 \\
&\quad - \frac{s_1}{s_k} \sum_{y_0} p(y_0|y_k) \sum_x p(x|y_k, y_0) \log \frac{p(x|y_k, y_0)}{q(x|u_k, y_0)} \frac{1}{p(x|y_k, y_0)} - \sum_{y_0} p(y_0|y_k) \log \frac{p(y_0|y_k)}{q(y_0|u_k)} \frac{1}{p(y_0|y_k)} \\
&= \log q(u_k) - \psi_k(u_k, y_k) + \tilde{\lambda}_k(y_k), \tag{68}
\end{aligned}$$

where  $\psi_k(u_k, y_k)$ ,  $k = 1, 2$ , are given by (35), and  $\tilde{\lambda}_k(y_k)$  contains all terms independent of  $u_k$  for  $k = 1, 2$ . Then, we proceeded by rearranging (68) as follows

$$p(u_k|y_k) = e^{\tilde{\lambda}_k(y_k)} q(u_k) e^{-\psi_k(u_k, y_k)}, \quad \text{for } k = 1, 2. \tag{69}$$

Finally, the Lagrange multipliers  $\lambda_k(y_k)$  satisfying the KKT conditions are obtained by finding  $\tilde{\lambda}_k(y_k)$  such that  $\sum_{u_k} p(u_k|y_k) = 1$ ,  $k = 1, 2$ . Substituting in (69),  $p(u_k|y_k)$  can be found as in (34).

## APPENDIX VII

### DERIVATION OF THE UPDATE RULES OF ALGORITHM 2

In this section, we derive the update rules in Algorithm 2 and show that the Gaussian distribution is invariant to the update rules in Algorithm 1, in line with Theorem 2.

First, we recall that if  $(\mathbf{X}_1, \mathbf{X}_2)$  are jointly Gaussian, then

$$P_{\mathbf{X}_2|\mathbf{X}_1} \sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{x}_2|\mathbf{x}_1}, \boldsymbol{\Sigma}_{\mathbf{x}_2|\mathbf{x}_1}),$$

where  $\boldsymbol{\mu}_{\mathbf{x}_2|\mathbf{x}_1} := \mathbf{K}_{\mathbf{x}_2|\mathbf{x}_1} \mathbf{x}_1$ ,  $\mathbf{K}_{\mathbf{x}_2|\mathbf{x}_1} := \boldsymbol{\Sigma}_{\mathbf{x}_2, \mathbf{x}_1} \boldsymbol{\Sigma}_{\mathbf{x}_1}^{-1}$ .

Then, for  $\mathbf{Q}^{(t+1)}$  computed as in (32) from  $\mathbf{P}^{(t)}$ , which is a set of Gaussian distributions, we have

$$\begin{aligned}
Q_{\mathbf{X}|\mathbf{U}_1, \mathbf{U}_2, \mathbf{Y}_0} &\sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{u}_1, \mathbf{u}_2, \mathbf{y}_0}, \boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{u}_1, \mathbf{u}_2, \mathbf{y}_0}), & Q_{\mathbf{X}|\mathbf{U}_k, \mathbf{Y}_0} &\sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{u}_k, \mathbf{y}_0}, \boldsymbol{\Sigma}_{\mathbf{x}|\mathbf{u}_k, \mathbf{y}_0}), \\
Q_{\mathbf{Y}_0|\mathbf{U}_k} &\sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{y}_0|\mathbf{u}_k}, \boldsymbol{\Sigma}_{\mathbf{y}_0|\mathbf{u}_k}), & Q_{\mathbf{U}_k} &\sim \mathcal{CN}(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{u}_k}).
\end{aligned}$$

Next, we look at the update  $\mathbf{P}^{(t+1)}$  as in (34) from given  $\mathbf{Q}^{(t+1)}$ . To compute  $\psi_k(\mathbf{u}_k^t, \mathbf{y}_k)$ , first, we note that

$$\mathbb{E}_{\mathbf{U}_{\bar{k}}, \mathbf{Y}_0 | \mathbf{y}_k} D(P_{\mathbf{X} | \mathbf{y}_k, \mathbf{U}_{\bar{k}}, \mathbf{Y}_0} \| Q_{\mathbf{X} | \mathbf{u}_k, \mathbf{U}_{\bar{k}}, \mathbf{Y}_0}) = D(P_{\mathbf{U}_{\bar{k}}, \mathbf{X}, \mathbf{Y}_0 | \mathbf{y}_k} \| Q_{\mathbf{U}_{\bar{k}}, \mathbf{X}, \mathbf{Y}_0 | \mathbf{u}_k}) - D(P_{\mathbf{U}_{\bar{k}}, \mathbf{Y}_0 | \mathbf{y}_k} \| Q_{\mathbf{U}_{\bar{k}}, \mathbf{Y}_0 | \mathbf{u}_k}), \quad (70)$$

$$\mathbb{E}_{\mathbf{Y}_0 | \mathbf{y}_k} D(P_{\mathbf{X} | \mathbf{y}_k, \mathbf{Y}_0} \| Q_{\mathbf{X} | \mathbf{u}_k, \mathbf{Y}_0}) = D(P_{\mathbf{X}, \mathbf{Y}_0 | \mathbf{y}_k} \| Q_{\mathbf{X}, \mathbf{Y}_0 | \mathbf{u}_k}) - D(P_{\mathbf{Y}_0 | \mathbf{y}_k} \| Q_{\mathbf{Y}_0 | \mathbf{u}_k}),$$

and that for two multivariate Gaussian distributions, i.e.,  $P_{\mathbf{X}_1} \sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{x}_1}, \boldsymbol{\Sigma}_{\mathbf{x}_1})$  and  $P_{\mathbf{X}_2} \sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{x}_2}, \boldsymbol{\Sigma}_{\mathbf{x}_2})$  in  $\mathbb{C}^N$ ,

$$D(P_{\mathbf{X}_1} \| P_{\mathbf{X}_2}) = (\boldsymbol{\mu}_{\mathbf{x}_1} - \boldsymbol{\mu}_{\mathbf{x}_2})^\dagger \boldsymbol{\Sigma}_{\mathbf{x}_2}^{-1} (\boldsymbol{\mu}_{\mathbf{x}_1} - \boldsymbol{\mu}_{\mathbf{x}_2}) + \log |\boldsymbol{\Sigma}_{\mathbf{x}_2} \boldsymbol{\Sigma}_{\mathbf{x}_1}^{-1}| + \text{tr}(\boldsymbol{\Sigma}_{\mathbf{x}_2}^{-1} \boldsymbol{\Sigma}_{\mathbf{x}_1}) - N. \quad (71)$$

Applying (70) and (71) in (35) and noting that all involved distributions are Gaussian, it follows that  $\psi_k(\mathbf{u}_k^t, \mathbf{y}_k)$  is a quadratic form. Then, since  $q^{(t)}(\mathbf{u}_k)$  is also Gaussian, the product  $\log(q^{(t)}(\mathbf{u}_k) \exp(-\psi_k(\mathbf{u}_k^t, \mathbf{y}_k)))$  is also a quadratic form, and identifying constant, first and second order terms, we can write

$$\log p^{(t+1)}(\mathbf{u}_k | \mathbf{y}_k) = -(\mathbf{u}_k - \boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k})^\dagger \boldsymbol{\Sigma}_{\mathbf{z}_k^{t+1}}^{-1} (\mathbf{u}_k - \boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k}) + Z(\mathbf{y}_k),$$

where

$$\begin{aligned} \boldsymbol{\Sigma}_{\mathbf{z}_k^{t+1}}^{-1} &= \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} + \frac{1-s_1}{s_k} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t} - \frac{1-s_1}{s_k} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{u}_k^t} \\ &\quad + \frac{s_1}{s_k} \mathbf{K}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t} + \frac{s_k - s_1}{s_k} \mathbf{K}_{\mathbf{y}_0 | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{\mathbf{y}_0 | \mathbf{u}_k^t}^{-1} \mathbf{K}_{\mathbf{y}_0 | \mathbf{u}_k^t} \end{aligned} \quad (72)$$

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k} &= \boldsymbol{\Sigma}_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0) | \mathbf{y}_k} - \frac{1-s_1}{s_k} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{u}_k^t, \mathbf{y}_0) | \mathbf{y}_k} \right. \\ &\quad \left. + \frac{s_1}{s_k} \mathbf{K}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{u}_k^t}^{-1} \mathbf{K}_{(\mathbf{x}, \mathbf{y}_0) | \mathbf{y}_k} + \frac{s_k - s_1}{s_k} \mathbf{K}_{\mathbf{y}_0 | \mathbf{u}_k^t}^\dagger \boldsymbol{\Sigma}_{\mathbf{y}_0 | \mathbf{u}_k^t}^{-1} \mathbf{K}_{\mathbf{y}_0 | \mathbf{y}_k} \right) \mathbf{y}_k. \end{aligned} \quad (73)$$

This shows that  $p^{(t+1)}(\mathbf{u}_k | \mathbf{y}_k)$  is a Gaussian distribution and that  $\mathbf{U}_k^{t+1}$  is distributed as  $\mathbf{U}_k^{t+1} \sim \mathcal{CN}(\boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k}, \boldsymbol{\Sigma}_{\mathbf{z}_k^{t+1}})$ .

Next, we simplify (72) to obtain the update rule (38a). From the matrix inversion lemma, similarly to [26], for  $(\mathbf{X}_1, \mathbf{X}_2)$  jointly Gaussian we have

$$\boldsymbol{\Sigma}_{\mathbf{x}_2 | \mathbf{x}_1}^{-1} = \boldsymbol{\Sigma}_{\mathbf{x}_2}^{-1} + \mathbf{K}_{\mathbf{x}_1 | \mathbf{x}_2}^\dagger \boldsymbol{\Sigma}_{\mathbf{x}_1 | \mathbf{x}_2}^{-1} \mathbf{K}_{\mathbf{x}_1 | \mathbf{x}_2}. \quad (74)$$

Applying (74) in (72), we have

$$\begin{aligned} \boldsymbol{\Sigma}_{\mathbf{z}_k^{t+1}}^{-1} &= \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} + \frac{1-s_1}{s_k} \left( \boldsymbol{\Sigma}_{\mathbf{u}_k^t | (\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)}^{-1} - \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} \right) - \frac{1-s_1}{s_k} \left( \boldsymbol{\Sigma}_{\mathbf{u}_k^t | (\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} - \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} \right) \\ &\quad + \frac{s_1}{s_k} \left( \boldsymbol{\Sigma}_{\mathbf{u}_k^t | (\mathbf{x}, \mathbf{y}_0)}^{-1} - \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} \right) + \frac{s_k - s_1}{s_k} \left( \boldsymbol{\Sigma}_{\mathbf{u}_k^t | \mathbf{y}_0}^{-1} - \boldsymbol{\Sigma}_{\mathbf{u}_k^t}^{-1} \right) \\ &\stackrel{(a)}{=} \frac{1}{s_k} \boldsymbol{\Sigma}_{\mathbf{u}_k^t | (\mathbf{x}, \mathbf{y}_0)}^{-1} - \frac{1-s_1}{s_k} \boldsymbol{\Sigma}_{\mathbf{u}_k^t | (\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} + \frac{s_k - s_1}{s_k} \boldsymbol{\Sigma}_{\mathbf{u}_k^t | \mathbf{y}_0}^{-1}, \end{aligned}$$

where (a) is due to the Markov chain  $\mathbf{U}_1 \ominus \mathbf{X} \ominus \mathbf{U}_2$ . We obtain (38a) by taking the inverse of both sides of (a).

Also from the matrix inversion lemma [26], for  $(\mathbf{X}_1, \mathbf{X}_2)$  jointly Gaussian we have

$$\Sigma_{\mathbf{x}_1}^{-1} \Sigma_{\mathbf{x}_1, \mathbf{x}_2} \Sigma_{\mathbf{x}_2 | \mathbf{x}_1}^{-1} = \Sigma_{\mathbf{x}_1 | \mathbf{x}_2}^{-1} \Sigma_{\mathbf{x}_1, \mathbf{x}_2} \Sigma_{\mathbf{x}_2}^{-1}. \quad (75)$$

Now, we simplify (73) to obtain the update rule (38b) as follows

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k} &= \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, (\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)}^{-1} | \mathbf{u}_k^t \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right. \\ &\quad - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, (\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} | \mathbf{u}_k^t \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \\ &\quad \left. + \frac{s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, (\mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{x}, \mathbf{y}_0)}^{-1} | \mathbf{u}_k^t \Sigma_{(\mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, \mathbf{y}_0} \Sigma_{\mathbf{y}_0}^{-1} | \mathbf{u}_k^t \Sigma_{\mathbf{y}_0, \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right) \mathbf{y}_k \\ &\stackrel{(a)}{=} \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \Sigma_{\mathbf{u}_k^t, (\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right. \\ &\quad - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{\mathbf{u}_k^t, (\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \\ &\quad \left. + \frac{s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0)} \Sigma_{\mathbf{u}_k^t, (\mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{x}, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, \mathbf{y}_0} \Sigma_{\mathbf{y}_0}^{-1} \Sigma_{\mathbf{y}_0, \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right) \mathbf{y}_k \\ &\stackrel{(b)}{=} \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \mathbf{A}_k^t \Sigma_{\mathbf{y}_k, (\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right. \\ &\quad - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)} \mathbf{A}_k^t \Sigma_{\mathbf{y}_k, (\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \\ &\quad \left. + \frac{s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0)} \mathbf{A}_k^t \Sigma_{\mathbf{y}_k, (\mathbf{x}, \mathbf{y}_0)} \Sigma_{(\mathbf{x}, \mathbf{y}_0)}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0), \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, \mathbf{y}_0} \mathbf{A}_k^t \Sigma_{\mathbf{y}_k, \mathbf{y}_0} \Sigma_{\mathbf{y}_0}^{-1} \Sigma_{\mathbf{y}_0, \mathbf{y}_k} \Sigma_{\mathbf{y}_k}^{-1} \right) \mathbf{y}_k \\ &\stackrel{(c)}{=} \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)} \mathbf{A}_k^t (\Sigma_{\mathbf{y}_k} - \Sigma_{\mathbf{y}_k | (\mathbf{u}_k^t, \mathbf{x}, \mathbf{y}_0)}) \Sigma_{\mathbf{y}_k}^{-1} \right. \\ &\quad - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)} \mathbf{A}_k^t (\Sigma_{\mathbf{y}_k} - \Sigma_{\mathbf{y}_k | (\mathbf{u}_k^t, \mathbf{y}_0)}) \Sigma_{\mathbf{y}_k}^{-1} \\ &\quad \left. + \frac{s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0)} \mathbf{A}_k^t (\Sigma_{\mathbf{y}_k} - \Sigma_{\mathbf{y}_k | (\mathbf{x}, \mathbf{y}_0)}) \Sigma_{\mathbf{y}_k}^{-1} + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, \mathbf{y}_0} \mathbf{A}_k^t (\Sigma_{\mathbf{y}_k} - \Sigma_{\mathbf{y}_k | \mathbf{y}_0}) \Sigma_{\mathbf{y}_k}^{-1} \right) \mathbf{y}_k \\ &\stackrel{(d)}{=} \Sigma_{\mathbf{z}_k^{t+1}} \left( \frac{1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{x}, \mathbf{y}_0)} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | (\mathbf{x}, \mathbf{y}_0)} \Sigma_{\mathbf{y}_k}^{-1}) - \frac{1-s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{(\mathbf{u}_k^t, \mathbf{y}_0)} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | (\mathbf{u}_k^t, \mathbf{y}_0)} \Sigma_{\mathbf{y}_k}^{-1}) \right. \\ &\quad \left. + \frac{s_k - s_1}{s_k} \Sigma_{\mathbf{u}_k^t}^{-1} \Sigma_{\mathbf{u}_k^t, \mathbf{y}_0} \mathbf{A}_k^t (\mathbf{I} - \Sigma_{\mathbf{y}_k | \mathbf{y}_0} \Sigma_{\mathbf{y}_k}^{-1}) \right) \mathbf{y}_k, \end{aligned}$$

where (a) follows from (75); (b) follows from the relation  $\Sigma_{\mathbf{u}_k, \mathbf{y}_0} = \mathbf{A}_k \Sigma_{\mathbf{y}_k, \mathbf{y}_0}$ ; (c) is due the definition of  $\Sigma_{\mathbf{x}_1 | \mathbf{x}_2}$ ; and (d) is due to the Markov chain  $\mathbf{U}_1 \ominus \mathbf{X} \ominus \mathbf{U}_2$ . Equation (38b) follows by noting that  $\boldsymbol{\mu}_{\mathbf{u}_k^{t+1} | \mathbf{y}_k} = \mathbf{A}_k^{t+1} \mathbf{y}_k$ .

## REFERENCES

- [1] Y. Ugur, I. E. Aguerri, and A. Zaidi, "A generalization of Blahut-Arimoto algorithm to compute rate-distortion regions of multiterminal source coding under logarithmic loss," in *Proc. of IEEE Inf. Theory Workshop*, Nov. 2017, pp. 349–353.
- [2] —, "Vector Gaussian CEO problem under logarithmic loss," in *Proc. of IEEE Inf. Theory Workshop*, Nov. 2018, pp. 515–519.
- [3] Y. Oohama, "Rate-distortion theory for Gaussian multiterminal source coding systems with several side informations at the decoder," *IEEE Trans. Inf. Theory*, vol. 51, no. 7, pp. 2577–2593, Jul. 2005.
- [4] V. Prabhakaran, D. Tse, and K. Ramachandran, "Rate region of the quadratic Gaussian CEO problem," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun.-Jul. 2004, p. 117.
- [5] J. Chen and J. Wang, "On the vector Gaussian CEO problem," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jul.-Aug. 2011, pp. 2050–2054.
- [6] J. Wang and J. Chen, "On the vector Gaussian  $L$ -terminal CEO problem," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jul. 2012, pp. 571–575.
- [7] T. Liu and P. Viswanath, "An extremal inequality motivated by multiterminal information-theoretic problems," *IEEE Trans. Inf. Theory*, vol. 53, no. 5, pp. 1839–1851, May 2007.
- [8] Y. Xu and Q. Wang, "Rate region of the vector Gaussian CEO problem with the trace distortion constraint," *IEEE Trans. Inf. Theory*, vol. 62, no. 4, pp. 1823–1835, Apr. 2016.
- [9] T. A. Courtade and R. D. Wesel, "Multiterminal source coding with an entropy-based distortion measure," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jul.-Aug. 2011, pp. 2040–2044.
- [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] J. Jiao, T. A. Courtade, K. Venkat, and T. Weissman, "Justification of logarithmic loss via the benefit of side information," *IEEE Trans. Inf. Theory*, vol. 61, no. 10, pp. 5357–5365, Oct. 2015.
- [12] A. No and T. Weissman, "Universality of logarithmic loss in lossy compression," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2015, pp. 2166–2170.
- [13] Y. Shkel, M. Raginsky, and S. Verdú, "Universal lossy compression under logarithmic loss," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2017, pp. 1157–1161.
- [14] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," in *Proc. of the 37th Annu. Allerton Conf. Commun., Control and Comput.*, 1999, pp. 368–377.
- [15] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning and Games*. New York, USA: Cambridge Univ. Press, 2006.
- [16] T. Andre, M. Antonini, M. Barlaud, and R. M. Gray, "Entropy-based distortion measure for image coding," in *Proc. of IEEE Int. Conf. Image Process.*, Oct. 2006, pp. 1157–1160.
- [17] K. Kittichokechai, Y.-K. Chia, T. J. Oechtering, M. Skoglund, and T. Weissman, "Secure source coding with a public helper," *IEEE Trans. Inf. Theory*, vol. 62, no. 7, pp. 3930–3949, Jul. 2016.
- [18] E. Ekrem and S. Ulukus, "An outer bound for the vector Gaussian CEO problem," *IEEE Trans. Inf. Theory*, vol. 60, no. 11, pp. 6870–6887, Nov. 2014.
- [19] S. Tavildar and P. Viswanath, "On the sum-rate of the vector Gaussian CEO problem," in *Proc. of 39th Asilomar Conf. Signals, Syst. Comput.*, Oct.-Nov. 2005, pp. 3–7.
- [20] H. Weingarten, Y. Steinberg, and S. Shamai (Shitz), "The capacity region of the Gaussian multiple-input multiple-output broadcast channel," *IEEE Trans. Inf. Theory*, vol. 52, no. 9, pp. 3936–3964, Sep. 2006.
- [21] D. P. Palomar, J. M. Cioffi, and M. A. Lagunas, "Joint Tx-Rx beamforming design for multicarrier MIMO channels: A unified framework for convex optimization," *IEEE Trans. Signal Process.*, vol. 51, no. 9, pp. 2381–2401, Sep. 2003.
- [22] A. Scaglione, P. Stoica, S. Barbarossa, G. B. Giannakis, and H. Sampath, "Optimal designs for space-time linear precoders and decoders," *IEEE Trans. Signal Process.*, vol. 50, no. 5, pp. 1051–1064, May 2002.
- [23] P. Harremoës and N. Tishby, "The information bottleneck revisited or how to choose a good distortion measure," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2007, pp. 566–570.

- [24] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *IEEE Trans. Inf. Theory*, vol. IT-18, no. 4, pp. 460–473, Jul. 1972.
- [25] S. Arimoto, "An algorithm for computing the capacity of arbitrary discrete memoryless channels," *IEEE Trans. Inf. Theory*, vol. IT-18, no. 1, pp. 14–20, Jan. 1972.
- [26] G. Chechik, A. Globerson, N. Tishby, and Y. Weiss, "Information bottleneck for Gaussian variables," *J. of Mach. Learn. Research*, vol. 6, pp. 165–188, Jan. 2005.
- [27] A. Winkelbauer and G. Matz, "Rate-information-optimal Gaussian channel output compression," in *Proc. of the 48th Annu. Conf. Inf. Sciences and Sys.*, Aug. 2014.
- [28] A. Winkelbauer, S. Farthofer, and G. Matz, "The rate-information trade-off for Gaussian vector channels," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2014, pp. 2849–2853.
- [29] S. Cheng, V. Stankovic, and Z. Xiong, "Computing the channel capacity and rate-distortion function with two-sided state information," *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4418–4425, Dec. 2005.
- [30] M. Chiang and S. Boyd, "Geometric programming duals of channel capacity and rate distortion," *IEEE Trans. Inf. Theory*, vol. 50, no. 2, pp. 245–258, Feb. 2004.
- [31] F. Dupuis, W. Yu, and F. M. J. Willems, "Blahut-Arimoto algorithms for computing channel capacity and rate-distortion with side information," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun.-Jul. 2004, p. 181.
- [32] M. Rezaeian and A. Grant, "A generalization of Arimoto-Blahut algorithm," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun.-Jul. 2004, p. 180.
- [33] A. Painsky and G. Wornell, "On the universality of the logistic loss function," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2018, pp. 936–940.
- [34] C. T. Li, X. Wu, A. Ozgur, and A. E. Gamal, "Minimax learning for remote prediction," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2018, pp. 541–545.
- [35] C. Tian and J. Chen, "Remote vector Gaussian source coding with decoder side information under mutual information and distortion constraints," *IEEE Trans. Inf. Theory*, vol. 55, no. 10, pp. 4676–4680, Oct. 2009.
- [36] A. Sanderovich, S. Shamai (Shitz), Y. Steinberg, and G. Kramer, "Communication via decentralized processing," *IEEE Trans. Inf. Theory*, vol. 54, no. 7, pp. 3008–3023, Jul. 2008.
- [37] O. Simeone, E. Erkip, and S. Shamai (Shitz), "On codebook information for interference relay channels with out-of-band relaying," *IEEE Trans. Inf. Theory*, vol. 57, no. 5, pp. 2880–2888, May 2011.
- [38] I. E. Aguerri, A. Zaidi, G. Caire, and S. Shamai (Shitz), "On the capacity of cloud radio access networks with oblivious relaying," in *Proc. of IEEE Int. Symp. Inf. Theory*, Jun. 2017, pp. 2068–2072.
- [39] —, "On the capacity of cloud radio access networks with oblivious relaying," *IEEE Trans. Inf. Theory*, vol. 65, no. 7, pp. 4575–4596, July 2019.
- [40] F. P. Calmon, A. Makhdoumi, M. Medard, M. Varia, M. Christiansen, and K.-D. Duffy, "Principal inertia components and applications," *IEEE Trans. Inf. Theory*, vol. 63, no. 8, pp. 5011–5038, Jul. 2017.
- [41] R. Ahlswede and I. Csiszar, "Hypothesis testing with communication constraints," *IEEE Trans. Inf. Theory*, vol. IT-32, no. 4, pp. 533–542, Jul. 1986.
- [42] T. Han, "Hypothesis testing with multiterminal data compression," *IEEE Trans. Inf. Theory*, vol. IT-33, no. 6, pp. 759–772, Nov. 1987.
- [43] M. S. Rahman and A. B. Wagner, "On the optimality of binning for distributed hypothesis testing," *IEEE Trans. Inf. Theory*, vol. 58, no. 10, pp. 6282–6303, Oct. 2012.
- [44] C. Tian and J. Chen, "Successive refinement for hypothesis testing and lossless one-helper problem," *IEEE Trans. Inf. Theory*, vol. 54, no. 10, pp. 4666–4681, Oct. 2008.
- [45] S. Salehkalaibar, M. Wigger, and R. Timo, "On hypothesis testing against conditional independence with multiple decision centers," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2409–2420, Jun. 2018.

- [46] Y. Zhou, Y. Xu, W. Yu, and J. Chen, "On the optimal fronthaul compression and decoding strategies for uplink cloud radio access networks," *IEEE Trans. Inf. Theory*, vol. 62, no. 12, pp. 7402–7418, Dec. 2016.
- [47] T. A. Courtade, "Gaussian multiterminal source coding through the lens of logarithmic loss," in *Inf. Theory and Appl. Workshop*, 2015.
- [48] —, "A strong entropy power inequality," *IEEE Trans. Inf. Theory*, vol. 64, no. 4, pp. 2173–2192, Apr. 2018.
- [49] A. B. Wagner, S. Tavildar, and P. Viswanath, "Rate region of the quadratic Gaussian two-encoder source-coding problem," *IEEE Trans. Inf. Theory*, vol. 54, no. 5, pp. 1938–1961, May 2008.
- [50] T. A. Courtade and J. Jiao, "An extremal inequality for long Markov chains," in *Proc. of the 52nd Annu. Allerton Conf. Commun., Control and Comput.*, 2014.
- [51] N. Slonim and N. Tishby, "The power of word clusters for text classification," in *Proc. of 23rd European Colloq. Inf. Retrieval Research*, 2001, pp. 191–200.
- [52] Y. Baram, R. El-Yaniv, and K. Luz, "Online choice of active learning algorithms," *J. of Mach. Learn. Research*, vol. 5, pp. 255–291, Mar. 2004.
- [53] I. E. Aguerri and A. Zaidi, "Distributed information bottleneck method for discrete and Gaussian sources," in *Proc. of IEEE Int. Zurich Seminar Inf. and Commun.*, Feb. 2018.
- [54] —, "Distributed variational representation learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.
- [55] D. Russo and J. Zou, "How much does your data exploration overfit? Controlling bias via information usage," *arXiv: 1511.05219*, 2015.
- [56] A. Xu and M. Raginsky, "Information-theoretic analysis of generalization capability of learning algorithms," in *Proc. of Conf. Neural Inf. Process. Sys.*, 2017, pp. 2524–2533.
- [57] A. R. Asadi, E. Abbe, and S. Verdú, "Chaining mutual information and tightening generalization bounds," in *Proc. of the 32nd Conf. Neural Inf. Process. Syst.*, 2018.
- [58] J. Chen and T. Berger, "Successive Wyner-Ziv coding scheme and its application to the quadratic Gaussian CEO problem," *IEEE Trans. Inf. Theory*, vol. 54, no. 4, pp. 1586–1603, Apr. 2008.
- [59] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM J. Optim.*, vol. 23, no. 2, pp. 1126–1153, Jun. 2013.
- [60] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming," <http://cvxr.com/cvx>, Mar. 2014.
- [61] T. A. Courtade, J. Jiao, and T. Weissman, "On an extremal data processing inequality for long Markov chains," in *Proc. of IEEE Int. Zurich Seminar Inf. and Commun.*, Feb. 2014, pp. 33–36.
- [62] T. A. Courtade, "Information masking and amplification: The source coding setting," in *Proc. of IEEE Int. Symp. Inf. Theory*, 2012, pp. 189–193.
- [63] A. Dembo, T. M. Cover, and J. A. Thomas, "Information theoretic inequalities," *IEEE Trans. Inf. Theory*, vol. 37, no. 6, pp. 1501–1518, Nov. 1991.
- [64] D. P. Palomar and S. Verdú, "Gradient of mutual information in linear vector gaussian channels," *IEEE Trans. Inf. Theory*, vol. 52, no. 1, pp. 141–154, Jan. 2006.