

Deterministic Annealing-Based Optimization for Zero-Delay Source-Channel Coding in Networks

Mustafa Said Mehmetoglu, *Student Member, IEEE*, Emrah Akyol, *Member, IEEE*,
and Kenneth Rose, *Fellow, IEEE*

Abstract—This paper studies the problem of global optimization of zero-delay source-channel codes that map between the source space and the channel space, under a given transmission power constraint and for the mean-square-error distortion. Particularly, we focus on two well-known network settings: the Wyner-Ziv setting where only a decoder has access to side information and the distributed setting where independent encoders transmit over independent channels to a central decoder. Prior work derived the necessary conditions for optimality of the encoder and decoder mappings, along with a greedy optimization algorithm that imposes these conditions iteratively, in conjunction with the heuristic noisy channel relaxation method to mitigate poor local minima. While noisy channel relaxation is arguably effective in simple settings, it fails to provide accurate global optimization in more complicated settings considered in this paper. We propose a powerful nonconvex optimization method based on the concept of deterministic annealing—which is derived from information theoretic principles and was successfully employed in several problems including vector quantization, classification, and regression. We present comparative numerical results that show strict superiority of the proposed method over greedy optimization methods as well as prior approaches in literature.

Index Terms—Joint source channel coding, deterministic annealing, estimation, distributed coding.

I. INTRODUCTION

WHILE IT is well known that finite-delay coding schemes do not achieve the asymptotic bounds in general (see, e.g., [1, Theorem 21] or [2]), the problem of obtaining the optimal coding schemes for finite delay is an important open problem with considerable practical implications [3]–[9]. Recently, there has been growing interest in utilizing zero-delay mappings in network applications, see, e.g., [10], [11] for coding over multiple access channels, [12]–[14] for distributed coding of correlated sources and [15], [16] for analog multiple description coding.

Manuscript received January 14, 2015; revised June 4, 2015, August 17, 2015, and September 16, 2015; accepted October 6, 2015. Date of publication October 26, 2015; date of current version December 15, 2015. This work was supported by the National Science Foundation under the grants CCF-1016861, CCF-1118075 and CCF-1320599. The material in this paper was presented in part at the IEEE Information Theory Workshop, Sevilla, Spain, September 2013, and the IEEE International Conference on Acoustics, Speech, and Signal Processing, Florence, Italy, May 2014. The associate editor coordinating the review of this paper and approving it for publication was V. Stankovic.

M. S. Mehmetoglu and K. Rose are with the Department of Electrical and Computer Engineering, University of California, Santa Barbara, Santa Barbara, CA 93106 USA (e-mail: mehmetoglu@ece.ucsb.edu; rose@ece.ucsb.edu).

E. Akyol is with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA (e-mail: akyol@illinois.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCOMM.2015.2494004

Until recently, there have been two main approaches to numerical optimization of the mappings: i) Optimization of the parameter set of a structured mapping [8], [9], [17], [18]. The performance of this approach is limited to the parametric form (structure) assumed. For example, in [19] saw-tooth like structure is assumed for the mapping in the Wyner-Ziv setting and parameters of such mapping are optimized. ii) Design based on power constrained channel optimized vector quantization where a discretized version of the problem is tackled using tools developed for vector quantization [5], [20], [21].

Our approach builds on recent prior work in our lab [22] where the problem is studied in the original analog (functional) domain, i.e., without discretization in the problem formulation and without any assumption of a parametrized mapping. In [22], necessary conditions for optimality of mappings were derived, noting that while such conditions have theoretical value, they generally identify local optima. They are practically useless in the case of highly complex cost surfaces. In other words, simple greedy methods that are based on iterative imposition of necessary conditions of optimality tend to get trapped in local minima. In [22], “noisy channel relaxation” (NCR) [23] was employed to mitigate this problem. As we show in this work, while NCR is rather sufficient for simple settings, using more advanced non-convex optimization tools improves the performance significantly in sophisticated network scenarios.

In this paper, we propose a method based on a powerful non-convex optimization framework, *deterministic annealing*, to numerically approach globally optimal zero-delay mappings in network scenarios. Our preliminary results appeared in [24], [25]. We particularly focus on scenarios given in Figure 1: The first case is a point-to-point source-channel coding with decoder side information (i.e., the decoder has access to side information that is correlated with the source). The second setting involves distributed (separate) coding and transmission of two correlated sources to a central decoder that reconstructs individual sources. We also consider the function computation problem, where the decoder estimates a function of the sources. This is of interest for certain applications such as a wireless sensor network deployed in order to compute a function of the measurements [26]–[30].

Deterministic annealing (DA) is derived within a probabilistic framework where the main idea is to introduce controlled randomization into the optimization process, yet deterministically optimize the appropriate expectation functionals. The application-specific cost is minimized at successive stages of decreasing randomness and a nonrandom solution is obtained

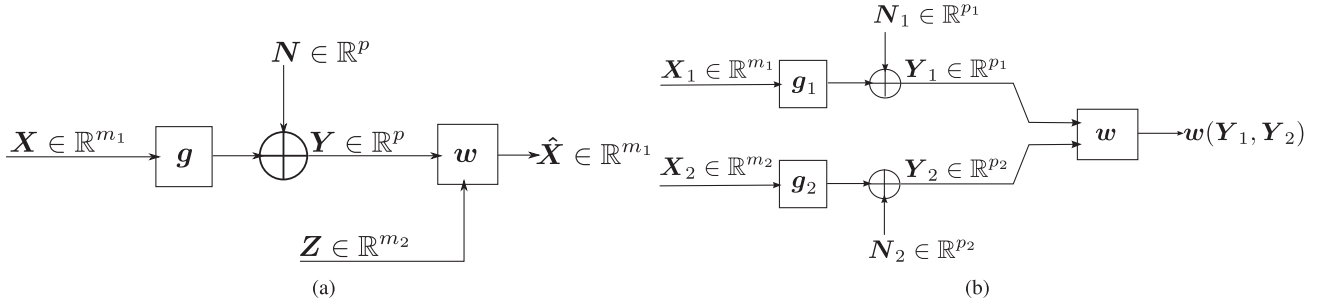


Fig. 1. Problem settings that we consider. (a) Decoder with side information. (b) Distributed coding setting.

while avoiding many poor local minima. Based on information theoretic principles, and motivated by analogies to statistical physics, DA has been successfully used in non-convex optimization problems including clustering [31], vector quantization [32], regression [33] and more (see review in [34]). We note that DA has been traditionally used in discrete settings, such as quantizer optimization, and integrating DA within the analog framework in here poses a significant challenge. There are many important advantages of the DA-based proposed method compared to prior work, including ability to avoid poor local minima and independence from initialization; and optimization in the original (analog) domain without any discretization or simplifying assumptions. Our approach improves significantly over prior approaches, some of which are NCR based [21], [22].

Having a powerful optimization method at hand, we analyze the structure of experimentally obtained mappings and investigate some conjectures made in prior work. For instance, one such conjecture was concerning the structure of optimal mappings in the side information setting, for which our results provide contradictory experimental evidence. Several practically important observations are made regarding the functional properties of the optimal mappings in network settings (see [35] for formal discussions of such properties in the point-to-point setting).

The rest of this paper is organized as follows. In Section II, we present preliminaries and the problem definition. In Section III and IV, we describe the proposed method. Experimental results are presented in Section V and concluding remarks are in Section VI.

II. PRELIMINARIES AND PROBLEM DEFINITIONS

A. Notations

Let \mathbb{R} , \mathbb{N} , and \mathbb{R}^+ denote the respective sets of real numbers, natural numbers, and positive real numbers. We represent scalars and random variables with lowercase and uppercase letters (e.g., x and X), column vectors and random column vectors with boldface lowercase and uppercase letters (e.g., \mathbf{x} and \mathbf{X}), respectively. $\|\cdot\|$ denotes L_2 norm operator. Let $\mathbb{E}(\cdot)$ and $\mathbb{P}(\cdot)$ denote the expectation and probability operators, respectively. The probability density function of the random variable X is $f_X(x)$. Let ∇ and ∇_x denote the gradient and partial gradient with respect to x , respectively. Let $f'(x) = \frac{df(x)}{dx}$ denote the first-order derivative of the continuously differentiable function

f . The Gaussian density with mean $\boldsymbol{\mu}$ and covariance matrix R is denoted as $\mathcal{N}(\boldsymbol{\mu}, R)$. We use natural logarithms which, in general, may be complex, and the integrals are, in general, Lebesgue integrals.

B. Problem Definition: Side Information

In the side information setting, given in Figure 1a, side information $\mathbf{Z} \in \mathbb{R}^{m_2}$ is available to the decoder, while source $\mathbf{X} \in \mathbb{R}^{m_1}$ is mapped to a channel input by the encoding function $\mathbf{g} : \mathbb{R}^{m_1} \rightarrow \mathbb{R}^p$ and transmitted over the channel with additive noise $\mathbf{N} \in \mathbb{R}^p$. The received channel output $\mathbf{Y} = \mathbf{g}(\mathbf{X}) + \mathbf{N}$ and side information \mathbf{Z} are mapped to the estimate $\hat{\mathbf{X}}$ by the decoding function $\mathbf{w} : \mathbb{R}^p \times \mathbb{R}^{m_2} \rightarrow \mathbb{R}^{m_1}$. The problem is to find optimal mappings \mathbf{g}, \mathbf{w} , where optimality is in the sense that they minimize MSE

$$D(\mathbf{g}, \mathbf{w}) = \mathbb{E}\{\|\mathbf{X} - \hat{\mathbf{X}}\|^2\}, \quad (1)$$

subject to some power constraint on the encoder

$$P(\mathbf{g}) = \mathbb{E}\{\|\mathbf{g}(\mathbf{X})\|^2\} \leq P_E \quad (2)$$

where $P_E > 0$ is the specified encoder power level. Simple time-sharing arguments show that D is a convex functional of P , hence the solution is achieved at $P = P_E$ (see [35] for details.) Converting to Lagrangian formulation, we define the following cost to be minimized

$$J = D(\mathbf{g}, \mathbf{w}) + \lambda(P(\mathbf{g}) - P_E) \quad (3)$$

where λ is a Lagrange multiplier corresponding to the power constraint on the encoder (we suppressed the dependence of J on \mathbf{g} and \mathbf{w}).

C. Problem Definition: Distributed Coding

The distributed coding setting, given in Figure 1b, has two sources $\mathbf{X}_1 \in \mathbb{R}^{m_1}$ and $\mathbf{X}_2 \in \mathbb{R}^{m_2}$ mapped to some channel input by the encoding functions $\mathbf{g}_i : \mathbb{R}^{m_i} \rightarrow \mathbb{R}^{p_i}$, and the decoder receives $\mathbf{Y}_i = \mathbf{g}_i(\mathbf{X}_i) + \mathbf{N}_i$ for $i = 1, 2$. In general, the decoder might have two type of objectives. In the first one, the decoder aims to reconstruct each source with minimum distortion. The decoder is defined as $\mathbf{w} : \mathbb{R}^{p_1} \times \mathbb{R}^{p_2} \rightarrow \mathbb{R}^{m_1} \times \mathbb{R}^{m_2}$ as it maps the received channel outputs to the estimates $\hat{\mathbf{X}}_i$ for $i = 1, 2$. For this case, we define distortion as

$$D(\mathbf{g}_1, \mathbf{g}_2, \mathbf{w}) = \mathbb{E}\{\|\mathbf{X}_1 - \hat{\mathbf{X}}_1\|^2 + \eta\|\mathbf{X}_2 - \hat{\mathbf{X}}_2\|^2\} \quad (4)$$

where $\eta \in \mathbb{R}^+$ is a given weight coefficient. The second type of problems involve function computation. Denoting the desired function as $\boldsymbol{y}(\boldsymbol{X}_1, \boldsymbol{X}_2) : \mathbb{R}^{m_1} \times \mathbb{R}^{m_2} \rightarrow \mathbb{R}^r$, the decoder is defined as $\boldsymbol{w} : \mathbb{R}^{p_1} \times \mathbb{R}^{p_2} \rightarrow \mathbb{R}^r$ and the cost is given by

$$D(\boldsymbol{g}_1, \boldsymbol{g}_2, \boldsymbol{w}) = \mathbb{E}\{\|\boldsymbol{y}(\boldsymbol{X}_1, \boldsymbol{X}_2) - \boldsymbol{w}(\boldsymbol{Y}_1, \boldsymbol{Y}_2)\|^2\}. \quad (5)$$

The problem, for both cases, is to find the mappings $\boldsymbol{g}_1, \boldsymbol{g}_2, \boldsymbol{w}$ that minimize the overall distortion (which is given in (4) or (5) depending on the objective) subject to power constraints on the encoders, which can be in two forms: Individual power constraints given by

$$P(\boldsymbol{g}_i) = \mathbb{E}\{\|\boldsymbol{g}_i(\boldsymbol{X}_i)\|^2\} \leq P_{T,i} \text{ for } i = 1, 2. \quad (6)$$

or a total power allocation to the encoders

$$\sum_{i=1}^2 P(\boldsymbol{g}_i) \leq P_T, \quad (7)$$

which offers the additional degree of freedom of optimizing power allocation to the encoders. For optimization purposes, we similarly define the following Lagrangian functional as the objective cost to be minimized

$$J = D + \sum_{i=1}^2 \lambda_i (P(\boldsymbol{g}_i) - P_{T,i}), \quad (8)$$

where $\lambda_i \in \mathbb{R}^+$, $i = 1, 2$, are Lagrange multipliers to impose the individual power constraints on the encoders in the first case. The total power constraint case corresponds to the special case of (8) with $\lambda_1 = \lambda_2 = \lambda$, i.e., the Lagrangian cost to minimize is

$$J = D + \lambda(P(\boldsymbol{g}_1) + P(\boldsymbol{g}_2) - P_T), \quad (9)$$

where λ controls the total power.

D. Prior Work: Necessary Conditions of Optimality and Greedy Descent Algorithms

Here, we summarize the relevant contributions of prior work (see [22] for more details). For the side information setting, let the encoder \boldsymbol{g} be fixed. Then, the optimal decoder is the MSE estimator of \boldsymbol{X} given $\boldsymbol{Z} = \boldsymbol{z}$ and $\boldsymbol{Y} = \boldsymbol{y}$:

$$\boldsymbol{w}(\boldsymbol{y}, \boldsymbol{z}) = \mathbb{E}\{\boldsymbol{X}|\boldsymbol{y}, \boldsymbol{z}\}. \quad (10)$$

Expanding the expressions for expectation and applying Bayes' rule, the optimal decoder can be written in terms of known quantities as

$$\boldsymbol{w}(\boldsymbol{y}, \boldsymbol{z}) = \frac{\int \boldsymbol{x} f_{\boldsymbol{X}, \boldsymbol{Z}}(\boldsymbol{x}, \boldsymbol{z}) f_{\boldsymbol{N}}(\boldsymbol{y} - \boldsymbol{g}(\boldsymbol{x})) d\boldsymbol{x}}{\int f_{\boldsymbol{X}, \boldsymbol{Z}}(\boldsymbol{x}, \boldsymbol{z}) f_{\boldsymbol{N}}(\boldsymbol{y} - \boldsymbol{g}(\boldsymbol{x})) d\boldsymbol{x}}, \quad (11)$$

where we used the fact that $f_{\boldsymbol{Y}|\boldsymbol{X}}(\boldsymbol{y}, \boldsymbol{x}) = f_{\boldsymbol{N}}(\boldsymbol{y} - \boldsymbol{g}(\boldsymbol{x}))$. For optimality of \boldsymbol{g} , assuming the decoder \boldsymbol{w} is fixed, a necessary condition is

$$\nabla_{\boldsymbol{g}} J(\boldsymbol{g}, \boldsymbol{w}) = 0, \quad (12)$$

where

$$\begin{aligned} \nabla_{\boldsymbol{g}} J(\boldsymbol{g}, \boldsymbol{w}) &= \lambda f_{\boldsymbol{X}}(\boldsymbol{x}) \boldsymbol{g}(\boldsymbol{x}) \\ &\quad - \mathbb{E}\{\boldsymbol{w}'(\boldsymbol{g}(\boldsymbol{x}) + \boldsymbol{N}, \boldsymbol{Z})(\boldsymbol{x} - \boldsymbol{w}(\boldsymbol{g}(\boldsymbol{x}) + \boldsymbol{N}, \boldsymbol{Z}))\}, \end{aligned} \quad (13)$$

and \boldsymbol{w}' denotes the Jacobian of \boldsymbol{w} with respect to its first argument (see [22] for proof).

Remark 1: Note that in the case of jointly Gaussian sources and Gaussian channel(s) with matched source-channel dimensions, linear mappings satisfy the necessary conditions of optimality, however, they are highly suboptimal, see, e.g., [22]. As we will see, careful optimization obtains considerably better mappings that are far from linear.

The necessary conditions of optimality for the distributed coding setting can be derived similarly, and are omitted for brevity, see [22]. Iteratively alternating between the imposition of individual necessary conditions of optimality will successively decrease the Lagrangian cost until a stationary point is reached. We refer to this method as ‘‘greedy descent’’. There is no reason to expect that a greedy descent algorithm will converge to the globally optimal solution. In fact, experiments show severe issues of local optima and strong dependence on initialization of such methods. As a remedy, the noisy channel relaxation (NCR) method of [23] was embedded in the algorithm in [22], i.e., the descent method was run at gradually decreasing levels of λ , wherein the result at each level serves as initialization for the next level of λ (see [23] for details). While such simple relaxations are effective in simple communication settings, the networked problems we consider here require a stronger optimization approach.

E. Asymptotically Achievable Limits

It is insightful to consider asymptotic bounds, which are obtained at infinite delay, while keeping in mind that the problem we consider is delay limited. Let $R(D)$ and $C(P)$ denote the source rate-distortion function and channel capacity, respectively. According to Shannon's source and channel coding theorems, the source can be compressed to $R(D)$ bits (per source sample) at distortion level D , and that $C(P)$ bits can be transmitted over the channel (per channel use) with arbitrarily low probability of error (see, e.g., [36]). The optimal coding scheme is the tandem combination of the optimal source and channel coding schemes, hence, by setting

$$R(D) = C(P), \quad (14)$$

one obtains a lower bound on the distortion of any source-channel coding scheme. For simplicity, we derive the expressions for the ‘‘optimum performance theoretically attainable’’ (OPTA) for Gaussian scalar source and noise. The channel capacity with additive white Gaussian noise is given by

$$C(P) = \frac{1}{2} \log \left(1 + \frac{P}{\sigma_N^2} \right), \quad (15)$$

where P is the transmission power and σ_N^2 is the noise variance.

For source-channel coding with decoder side information, OPTA can be obtained by equating Wyner-Ziv rate distortion function [37] to the channel capacity. The Wyner-Ziv rate distortion function of X , when Z serves as side information, and $(X, Z) \sim \mathcal{N}(\mathbf{0}, R_{X,Z})$ where $R_{X,Z} = \sigma_X^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ and σ_X^2, ρ are the variance and correlation coefficient, respectively, with $|\rho| \leq 1$ is:

$$R(D) = \max \left(0, \frac{1}{2} \log \frac{(1 - \rho^2)\sigma_X^2}{D} \right). \quad (16)$$

We plug (16) and (15) in (14) to obtain

$$D_{OPTA} = \frac{(1 - \rho^2)\sigma_X^2}{\left(1 + \frac{P_T}{\sigma_N^2}\right)}. \quad (17)$$

For quadratic Gaussian distributed source coding for sources $(X_1, X_2) \sim \mathcal{N}(\mathbf{0}, R_{X_1, X_2})$ where $R_{X_1, X_2} = \sigma_X^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ with $|\rho| \leq 1$, the complete rate distortion region satisfies the following inequalities [38]:

$$R_1 \geq \frac{1}{2} \log^+ \left(\frac{1 - \rho^2 + \rho^2 2^{-2R_2}}{D_1} \right) \quad (18)$$

$$R_2 \geq \frac{1}{2} \log^+ \left(\frac{1 - \rho^2 + \rho^2 2^{-2R_1}}{D_2} \right) \quad (19)$$

$$R_1 + R_2 \geq \frac{1}{2} \log^+ \left(\frac{(1 - \rho^2)\beta(D_1, D_2)}{2D_1 D_2} \right) \quad (20)$$

where $\log^+ x = \max(0, \log x)$ and

$$\beta(D_1, D_2) = 1 + \sqrt{1 + \frac{4\rho^2 D_1 D_2}{(1 - \rho^2)^2}}. \quad (21)$$

We set $R_i = C(P_i)$ for $i = 1, 2$, where $C(P)$ is given in (15) to obtain OPTA.

III. PROPOSED METHOD FOR SIDE INFORMATION SETTING

A. Overview

In this section, we develop the DA based method for the optimization of encoder and decoder mappings. Since the decoder is given in closed form, the method focuses on optimizing the encoder mapping. We first partition the input space of the encoder into partition cells and assign a local model to each of the cells. Next, the encoder output is made probabilistic by randomizing the partitions, i.e., input points are assigned to each local model according to some probability distribution. We then propose an optimization process where the (random) encoder is optimized (along with the decoder) while constraining the Shannon entropy. By gradually reducing the entropy to 0, we obtain the desired mappings.

B. Derivation of Proposed Method

We consider piecewise functions which approximate the desired mappings by partitioning the space and matching a sim-

ple local model to each region. Piecewise functions consist of two components: a space partition and a parametric local model per partition cell. First, the source space \mathbb{R}^m is partitioned into \mathcal{K} regions (cells) denoted \mathcal{R}_k^m . Each cell \mathcal{R}_k^m has an associated function \mathbf{g}_k which is parametrized (affine, lattice, etc.) and the parameter set is denoted by Λ_k . Thus, the encoding function can be written as

$$\mathbf{g}(\mathbf{x}) = \mathbf{g}_k(\mathbf{x}) \text{ for } \mathbf{x} \in \mathcal{R}_k^m \text{ and for } k = 1, \dots, \mathcal{K} \quad (22)$$

In (22), the selection of local model index k is deterministic for a given realization of \mathbf{X} , i.e., the output of the encoder only depends on \mathbf{X} . To derive a DA based approach, we introduce a random variable, K , that corresponds to random selection of index k . In other words, let the encoder randomly select the local model index k when it receives an input \mathbf{x} , according to the value of a random variable that we call K . For a given realization of \mathbf{X} , the output of the encoder is now given in probability as

$$\mathbf{g}(\mathbf{x}) = \mathbf{g}_k(\mathbf{x}) \text{ with probability } p_{K|\mathbf{X}}(k|\mathbf{x}). \quad (23)$$

The conditional probability $p_{K|\mathbf{X}}(k|\mathbf{x})$ is referred to as *association probability*, in the sense that it represents the probability of input point \mathbf{x} belonging to cell \mathcal{R}_k^m (thus, the source space partition is now random). The probability distribution that we introduce (and optimize) is $p_{K|\mathbf{X}}$ (not the joint $p_{\mathbf{X}, K}$) since the input distribution is given in the problem statement and is therefore fixed. The MSE cost and transmission power are still calculated as in (1) and (2), though the expectation is now taken over K in addition to what was done before. Let us now fix Λ_k and \mathbf{w} , and consider optimizing (3) with respect to $p_{K|\mathbf{X}}$. It is clear that the optimal $p_{K|\mathbf{X}}$ will implement 'hard' associations, that is, every point \mathbf{x} will be fully associated with the local model that makes the minimum contribution to cost¹. Although this is desirable eventually, in order to avoid poor local optima we impose and control the level of randomness, i.e., we introduce a constraint on the randomness of the encoder, which is measured by the Shannon entropy. The total entropy of the encoder is given by $H(\mathbf{X}, K) = H(\mathbf{X}) + H(K|\mathbf{X})$ and since $H(\mathbf{X})$ is constant (predetermined by the source), the entropic quantity of interest is the conditional entropy $H(K|\mathbf{X})$. This is also intuitively justified in the sense that the randomness we introduced into the problem is precisely captured by $p_{K|\mathbf{X}}$, hence can be measured and controlled by $H(K|\mathbf{X})$. We denote the randomness of the solution by H and define it as $H \triangleq H(K|\mathbf{X})$ where

$$H(K|\mathbf{X}) = -\mathbb{E}\{\log p_{K|\mathbf{X}}\}. \quad (24)$$

The problem is now recast as minimization of the expected cost with respect to parameters of local models, association

¹Therefore, the generalized search space of random encoders have the same global minimum as the original problem.

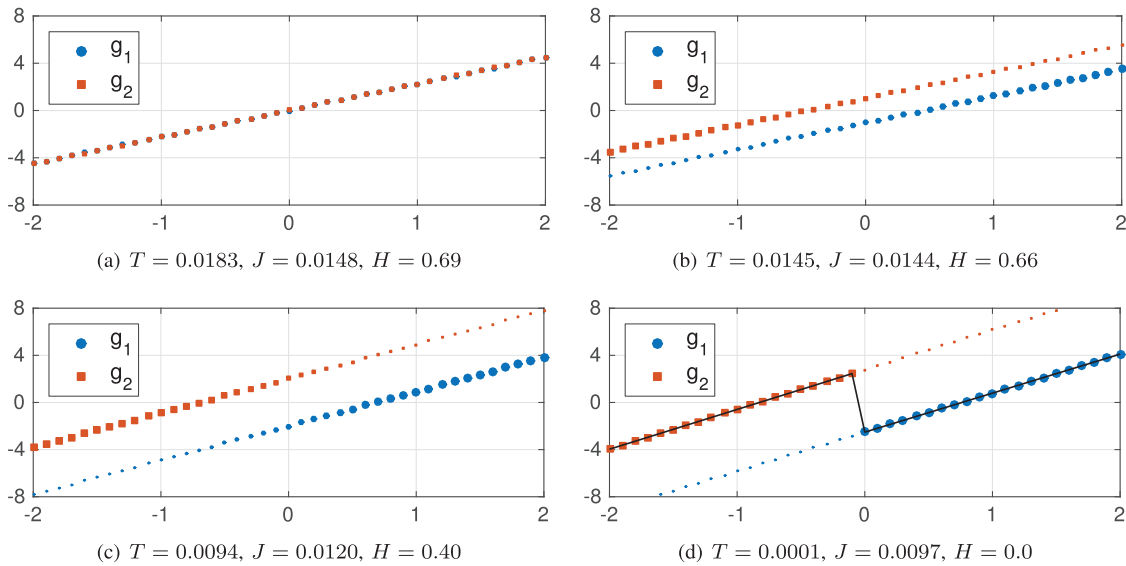


Fig. 2. The evolution of the encoder in the algorithm is demonstrated. The two models are shown by dotted lines and the sizes of dots are relative to the probability association at that input point. The line in (d) is the deterministic encoder obtained. $\mathcal{K} = 2$.

probabilities and decoder, subject to a constraint on the level of randomness of the system, i.e.,

$$\begin{aligned} & \underset{\Lambda_1, \dots, \Lambda_{\mathcal{K}}, p(1|\mathbf{x}), \dots, p(\mathcal{K}|\mathbf{x}), \mathbf{w}}{\text{minimize}} && J, \\ & \text{subject to} && H \geq H_0, \end{aligned}$$

where J is defined in (3) and H_0 specifies the minimum requirement on the entropy level. This constrained optimization problem can be reformulated by introducing Lagrange parameter $T \in \mathbb{R}^+$ to obtain the Lagrangian

$$F = J - T(H - H_0), \quad (25)$$

to be minimized. There are two important extremal points of this Lagrangian. First, for $T \rightarrow \infty$, the minimum F is obtained by maximizing the entropy, which is achieved by uniform association probabilities: $p_{\mathcal{K}|X}(k|\mathbf{x}) = 1/\mathcal{K}$ for all k and \mathbf{x} . Consequently, all local models equally account for all points and are identical once optimized, or effectively, there is a single *distinct* local model. Secondly, in the limit $T \rightarrow 0$, minimizing F corresponds to minimizing J directly, which produces a deterministic encoder. This intuitive observation can be verified by the expression for optimal $p_{\mathcal{K}|X}(k|\mathbf{x})$ given in Section III-D.

Although DA is derived from information theoretic principles, it is motivated by and has strong analogies to annealing processes in statistical physics (see [34] for details). We accordingly refer to the Lagrangian functional in (25) as (Helmholtz) free energy, and Lagrange parameter T as “temperature”.

C. Deterministic Annealing

The optimization method starts at a high value of T and gradually lowers it while minimizing F at each step. At high temperature, there is effectively a single distinct local model. As the temperature is decreased, a bifurcation point is reached where the current solution is no longer a minimum, so that

there exists a better solution with a higher number of distinct local models. Intuitively, at this temperature, the current solution is a saddle point where multiple local models are coincident (i.e., their parameters are same) and in order to move to a better solution, it is necessary to perturb the local models. Such bifurcations are referred to as “phase transitions” and the corresponding temperatures are called “critical temperatures”².

We present an example simulation in Figure 2 that illustrates the basics of the method, including phase transitions. Here the sources and channel are scalar, i.e., $m = n = 1$, g_k are selected as affine and $\mathcal{K} = 2$. When T is large, there is a single distinct local model. As we lower T , the system goes through a phase transition where the two local models split from each other (after a slight perturbation). The corresponding value of T is referred to as the first critical temperature. Note how entropy (H) is traded for reduction in cost (J).

Mappings with more than 2 local models can be obtained by starting with a larger \mathcal{K} . However, a computationally more efficient method that we employ here is as follows: We start with 1 local model and keep only the distinct local models, but duplicate and perturb them at each temperature. The duplicates will merge at every iteration until a critical temperature is reached, and will split into distinct models at a phase transition.

Although our method is derived in the general, continuous source and channel domain, in practical simulations we sample the source and noise distributions to allow numerical computation of integrals. The sampling is not “inherent” to the derived method and, in fact, can be adjusted during the algorithm run. We emphasize that this is in contrast with prior quantizer design based methods that are entirely formulated in a discrete setting.

The practical algorithm is initialized with a single local model. Since T must be set higher than the first critical temperature, we simply choose T large enough that during the first couple of temperatures, duplicated local models merge back,

²We omit the derivation of critical temperatures in this paper, see [34] for phase transition analysis in the general DA setting.

i.e., no phase transitions are observed. As the temperature is gradually lowered, we track the minimum, i.e., find the association probabilities $p_{K|X}(k|\mathbf{x})$, local model parameters Λ_k and decoder \mathbf{w} that minimize the Lagrangian F . As demonstrated, the system will go through phase transitions during which the number of local models, \mathcal{K} , increases. We stop when T is near 0 and perform “zero entropy iteration”, i.e., associate every source point with the “best” local model to obtain deterministic encoder. We accordingly give a brief sketch of the practical method in Algorithm 1. In Step 6, we employed an exponential cooling schedule. Update equations for Step 3 are given in the next section.

Algorithm 1 Proposed DA-Based Method

Inputs: Involved distributions, desired local model type, λ , α , ϵ , Δ_F , T_{min} , $\Delta_{\mathbf{g}}$.

Outputs: Optimized \mathbf{g} , \mathbf{w} .

Initialization: $T = T_{max}$, $\mathcal{K} = 1$, randomly chosen \mathbf{g}_1 . $J_{old} = J_{initial}$.

1. **Duplication:**

For each \mathbf{g}_i , create an identical local model \mathbf{g}_j .

$p(i|\mathbf{x}) \leftarrow \frac{p(i|\mathbf{x})}{2}$ and $p(j|\mathbf{x}) \leftarrow \frac{p(i|\mathbf{x})}{2}$.

$\mathcal{K} \leftarrow 2\mathcal{K}$.

2. **Perturbation:**

For each parameter $\phi_k \in \Lambda_k$, $\phi_k \leftarrow \phi_k + \epsilon R$, where R is standard Gaussian random variable.

3. **Thermal Equilibrium:**

Compute F and set $F_{old} \leftarrow F$.

3.1 Compute optimal \mathbf{w} using (30).

3.2 Compute optimal $p(k|\mathbf{x})$, $\forall k$ using (26).

3.3 Optimize Λ_k , $\forall k$ using (28).

3.4 Compute F . If $\frac{F - F_{old}}{F_{old}} \leq \Delta_F$, go to Step 4, otherwise $F_{old} \leftarrow F$ and go to Step 3.1.

4. **Model Size:**

If $d(\Lambda_i, \Lambda_j) < \Delta_{\mathbf{g}}$, where $d(\cdot, \cdot)$ is euclidean distance, remove \mathbf{g}_j and set $p(i|\mathbf{x}) \leftarrow p(i|\mathbf{x}) + p(j|\mathbf{x})$, $\forall i, j$.

$\mathcal{K} \leftarrow$ New model size.

5. **Stopping:**

Stop if $T \leq T_{min}$, otherwise go to Step 6.

6. **Cooling:**

$T \leftarrow T * \alpha$.

Go to Step 1.

D. Update Equations

The central part of the method is the minimization of free energy (F) by iteratively updating the association probabilities, local model parameters and decoders. The following theorem, whose proof is presented in the Appendix, states the update equations for association probabilities.

Theorem 1: At any temperature T , minimum free energy F is achieved when association probabilities are in the form of Gibbs distribution given as:

$$p(k|\mathbf{x}) = \frac{e^{-J_k(\mathbf{x})/T}}{\sum_{k'} e^{-J_{k'}(\mathbf{x})/T}} \quad \forall k, \quad (26)$$

where $J_k(\mathbf{x})$ is given by

$$J_k(\mathbf{x}) = \mathbb{E}\{\|\mathbf{x} - \mathbf{w}(\mathbf{g}_k(\mathbf{x}) + \mathbf{N}, \mathbf{Z})\|^2\} + \lambda \|\mathbf{g}_k(\mathbf{x})\|^2. \quad (27)$$

Remark 2: Theorem 1 is analogous to the principle of minimal free energy in statistical physics. A fundamental principle in statistical physics states that the minimum free energy is achieved when the system is at thermal equilibrium, at which point it is governed by Gibbs distribution.

The evolution of association probabilities, $p(k|\mathbf{x})$, during the annealing process can be observed from how (26) is changing with T . The following corollary confirms the intuitive explanation we provided earlier.

Corollary 1: As $T \rightarrow \infty$ (at a high temperature) the system is governed by uniform association probabilities and the entropy is maximum. As $T \rightarrow 0$, the associations become deterministic and the entropy is 0.

The optimal local model parameters cannot be obtained in closed form, hence we perform gradient descent search. A local model parameter $\phi_k \in \Lambda_k$ is updated according to

$$\phi_k \leftarrow \phi_k - \varphi \frac{\partial F}{\partial \phi_k} \quad (28)$$

where φ is selected by line search and the gradient can be obtained as

$$\frac{\partial F}{\partial \phi_k} = \frac{\partial J}{\partial \phi_k} = \int_{\mathbf{x}} f_X(\mathbf{x}) p(k|\mathbf{x}) \frac{\partial J_k(\mathbf{x})}{\partial \phi_k} d\mathbf{x}. \quad (29)$$

The derivative $\frac{\partial J_k(\mathbf{x})}{\partial \phi_k}$ is calculated numerically. The optimal decoder can be derived similar to (11):

$$\mathbf{w}(\mathbf{y}, \mathbf{z}) = \frac{\int_{\mathbf{x}} f_{X,Z}(\mathbf{x}, \mathbf{z}) \sum_k f_N(\mathbf{y} - \mathbf{g}_k(\mathbf{x})) p(k|\mathbf{x}) d\mathbf{x}}{\int_{\mathbf{x}} f_{X,Z}(\mathbf{x}, \mathbf{z}) \sum_k f_N(\mathbf{y} - \mathbf{g}_k(\mathbf{x})) p(k|\mathbf{x}) d\mathbf{x}}. \quad (30)$$

E. Design Complexity

Due to difficulties in estimating the time required for gradient descent, exact comparison of computational complexity of numerical optimization methods (including the method presented here and others referred to in Section II-D) is difficult and depends on the actual source-channel distributions as well as choice of various algorithm parameters. On the other hand, optimization of parametrized mappings (e.g., in [13]) is faster, but requires knowing the structure of a good solution, which can be obtained by methods such as the one presented here. In our experiments, the time required for DA was on the same order as that of NCR, albeit with a higher constant. Thus, better performance is obtained at the expense of slight increase in complexity.

IV. METHOD FOR DISTRIBUTED CODING

Although the method described in the previous section can be used for optimizing the distributed encoders separately (within separate annealing processes), we found that such a method

fails to avoid poor local minima as it fails to account for interaction between encoder optimizations. Instead, we develop a method here that optimizes the (random) encoders and decoders within a single annealing process. The resulting annealing method is a direct extension of the previous one, albeit with higher complexity due to the distributed nature of the problem.

We have two independent sets of partitions of input source space: \mathcal{K}_1 cells represented by $\mathcal{R}_{k_1}^m$ and \mathcal{K}_2 cells represented by $\mathcal{R}_{k_2}^m$. We define both encoders in this setting as

$$\mathbf{g}_i(\mathbf{x}_i) = \mathbf{g}_{k_i}(\mathbf{x}_i) \text{ for } \mathbf{x}_i \in \mathbb{R}_{k_i}^m, i = 1, 2. \quad (31)$$

Following the same procedure of randomization, we define random variables K_1 and K_2 along with association probabilities:

$$p(k_i|\mathbf{x}_i) \triangleq \mathbb{P}\{\mathbf{x}_i \in \mathbb{R}_{k_i}^m\}, \quad \forall k_i, \mathbf{x}_i, \text{ for } i = 1, 2. \quad (32)$$

The cost is to be minimized subject to the constraint on the joint entropy of the system. Noting that $K_1 \leftrightarrow X_1 \leftrightarrow X_2 \leftrightarrow K_2$ form a Markov chain by construction, we express the joint entropy as

$$H(X_1, K_1, X_2, K_2) = H(X_1, X_2) + H(K_1|X_1) + H(K_2|X_2). \quad (33)$$

Since $H(X_1, X_2)$ is a constant determined by the sources, we define $H \triangleq H(K_1|X_1) + H(K_2|X_2)$ where

$$H(K_i|X_i) = \mathbb{E}\{\log p(K_i|X_i)\} \text{ for } i = 1, 2, \quad (34)$$

and the free energy of the system is given by (25).

The algorithm sketch is similar to the side information setting and is not reproduced here. Since we optimize both encoders within the same annealing process, the same operations in the Algorithm are performed for both encoders, sequentially. The following theorem presents the optimal association probabilities for the distributed setting. The proof follows from the steps in the proof of Theorem 1 and omitted for brevity.

Theorem 2: At any temperature, minimum free energy (F) is achieved when the system is governed by Gibbs distribution given as:

$$p(k_i|\mathbf{x}_i) = \frac{e^{-J_{k_i}(\mathbf{x}_i)/T}}{\sum_{k'_i} e^{-J_{k'_i}(\mathbf{x}_i)/T}} \quad \text{for } i = 1, 2 \quad (35)$$

where

$$J_{k_i}(\mathbf{x}_i) = \mathbb{E}\{\|\mathbf{X}_1 - \hat{\mathbf{X}}_1\|^2 + \eta\|\mathbf{X}_2 - \hat{\mathbf{X}}_2\|^2 | \mathbf{X}_i = \mathbf{x}_i, K_i = k_i\} + \lambda_i \mathbf{g}_{k_i}^2(\mathbf{x}_i) \quad (36)$$

if the cost is defined as in (4), and

$$J_{k_i}(\mathbf{x}_i) = \mathbb{E}\{\|\boldsymbol{\gamma}(\mathbf{X}_1, \mathbf{X}_2) - \mathbf{w}(\mathbf{Y}_1, \mathbf{Y}_2)\|^2 | \mathbf{X}_i = \mathbf{x}_i, K_i = k_i\} + \lambda_i \mathbf{g}_{k_i}^2(\mathbf{x}_i) \quad (37)$$

if the cost is defined as in (5).

The parameters of local models can be optimized through gradient descent search. Optimal decoding is achieved similarly

as $\hat{\mathbf{X}}_i = \mathbb{E}\{\mathbf{X}_i | \mathbf{y}_1, \mathbf{y}_2\}$ for $i = 1, 2$ for first type of objective, and $\mathbf{w}(\mathbf{y}_1, \mathbf{y}_2) = \mathbb{E}\{\boldsymbol{\gamma}(\mathbf{X}_1, \mathbf{X}_2) | \mathbf{y}_1, \mathbf{y}_2\}$ for the second type. Both expressions can be written in terms of known quantities similar to that in (11).

V. EXPERIMENTAL RESULTS

While the proposed algorithm is general and directly applicable to any choice of source and channel dimensions, for conciseness of the results section, we assume that sources and channels are scalar. In this case, the encoder mapping is denoted as $g : \mathbb{R} \rightarrow \mathbb{R}$ and the local model functions g_k are selected as affine. In principle, the set of g_k can be chosen from any parametric model. Choosing a more complex model, such as a higher order polynomial, can potentially improve the performance of the algorithm, albeit with increased computational complexity. For the exponential cooling schedule, we set $\alpha = 0.95$, i.e., $T \leftarrow T * 0.95$. The performance of the proposed method is assessed by comparisons to the optimal affine solution, greedy method and NCR-based method developed in [22], as well as OPTA (for reference only, as OPTA requires infinite delay). For the NCR based method, we decrease λ (in distributed coding, we decrease λ_1 and λ_2 simultaneously) exponentially as $\lambda_{new} = \lambda_{old} * 0.8$ in 50 steps to the desired value.

The noise signals in all examples are chosen as independent zero-mean Gaussians with unit variance, i.e., $N \sim \mathcal{N}(0, 1)$, $N_1 \sim \mathcal{N}(0, 1)$, $N_2 \sim \mathcal{N}(0, 1)$. For numerical computations, we sample the source and noise distributions on a uniform grid with spacing $\Delta = 0.02$. We also impose bounded support (-5σ to $+5\sigma$), i.e., we neglect tails of infinite support distributions in the examples.

A. Side Information Setting

We first give examples for the Gaussian case, where the source and side information are jointly Gaussian, distributed according to $\mathcal{N}(\boldsymbol{\mu}, R)$ where $\boldsymbol{\mu} = [0, 0]$, $R = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$, and $|\rho| < 1$ is the correlation coefficient between source and side information. We define $\text{SNR} = 10 \log_{10}(1/D)$ and $\text{CSNR} = 10 \log_{10}(P(g))$.

Example mappings are given in Figure 3. We first note that the central characteristics observed in digital Wyner-Ziv mappings are captured by analog mappings as noted before (see, e.g., [21], [22]), in the sense of many-to-one mappings, where multiple source intervals are mapped to the same channel interval. We refer to each one-to-one section in these mappings as a ‘‘bin’’, in Figure 3a there are 5 bins in the interval shown (the meaning of bin here is different than in digital Wyner-Ziv mappings). The uncertainty about the source interval is resolved (significantly decreased) by the decoder using the side information. Since all variables are Gaussian and distortion measure is MSE, it is intuitively intriguing to investigate whether the optimal mappings have any parametric form or structure to be exploited in the design stage. For example, since in the absence of decoder side information optimal mappings are well known to be linear, one can expect to see linear mappings in each bin.

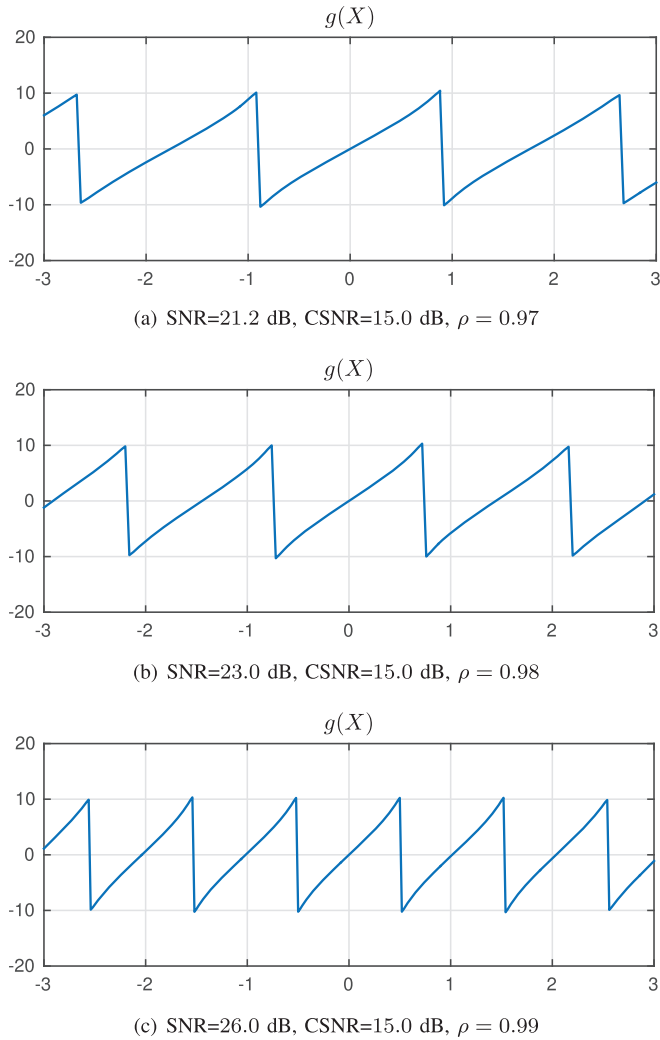


Fig. 3. Example encoder mappings, generated by DA, for the decoder side information setting, jointly Gaussian source and side information.

In fact, such parametric form was explicitly assumed in [19], and it was reported the optimized parametric mappings perform very close to the results obtained via NCR in [22]. Our numerical results demonstrate that each bin is non-linear as some nonlinearity can be observed especially near the ends of each bin, as opposed to the conjecture in [22].

From Figure 3 we see how the width of bins depends on the correlation between the source and side information. It can be seen that at higher correlation the bins are narrower. This is intuitively expected since, as the correlation increases, so does the benefit of side information in terms of distinguishing different bins. To exploit this capability, the encoder narrows the bins, which in turn reduces the power $\mathbb{E}\{g^2(X_1)\}$.

To illustrate the improvement of DA over NCR in the encoding mappings themselves, we present two mappings obtained by NCR in Figure 4. We emphasize that the performance of NCR depends on initial mappings, initial noise level and the noise-relaxation schedule. This dependence is illustrated in Figure 4, where in one case the shape of bins are different than those in DA and sub-optimal, and in the other the points of discontinuity are not optimal.

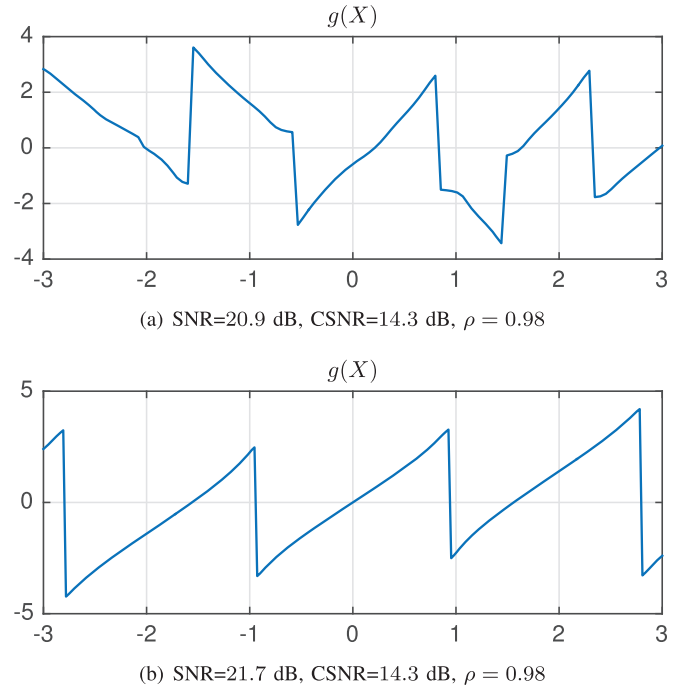


Fig. 4. Two results by NCR for side information setting. In (a) the bins do not have the optimal shape that was obtained by DA and in (b) the discontinuity points are not optimal.

We also give an example with a different source distribution, Gaussian mixture, in Figure 5:

$$(X_1, X_2) \sim \left(\frac{1}{2} \mathcal{N}(\boldsymbol{\mu}_1, R) + \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}_2, R) \right) \quad (38)$$

where $\boldsymbol{\mu}_1 = [-3, -3]$, $\boldsymbol{\mu}_2 = [3, 3]$ and $R = \begin{bmatrix} 1 & 0.95 \\ 0.95 & 1 \end{bmatrix}$. This distribution has two Gaussian “nodes” centered far from each other at $x = -3$ and $x = 3$. From an intuitive point of view, the optimum encoder can be viewed as two Wyner-Ziv like encoders, occupying the negative and positive halves of real line and both centered at the node centers. It is clear that for several source and channel distributions, optimal encoding mappings are many-to-one, i.e., this property is not unique to the Gaussian distribution.

The comparative performance results for different optimization techniques is given in Figure 6 for correlation coefficient $\rho = 0.99$. Since NCR performance depends on the initial conditions, we ran the NCR algorithm several times with different conditions and pick the mappings with best performance. Results from the greedy method are also presented in order to illustrate the abundance of locally optimum points and the difficulty of the optimization problem. Note that the proposed method is independent of the initialization and only run once. We also present the performance of OPTA as benchmark while noting that it is asymptotic and may require infinite delay. The performance of linear encoder and decoder is plotted as well, since it is also a local minimum (see Remark 1). It is important to note that the linear solution performs significantly worse than the non-linear mappings obtained.

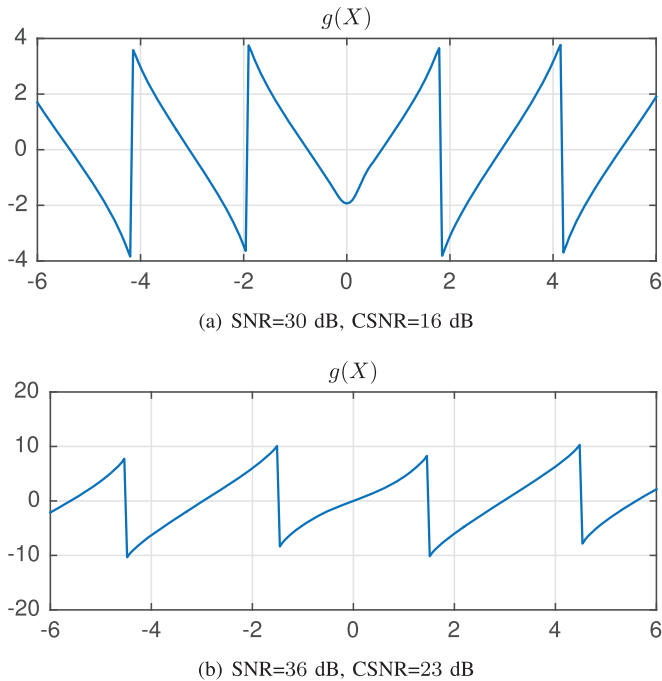


Fig. 5. Example encoder mappings, generated by DA, for Gaussian mixture distribution, side information setting.

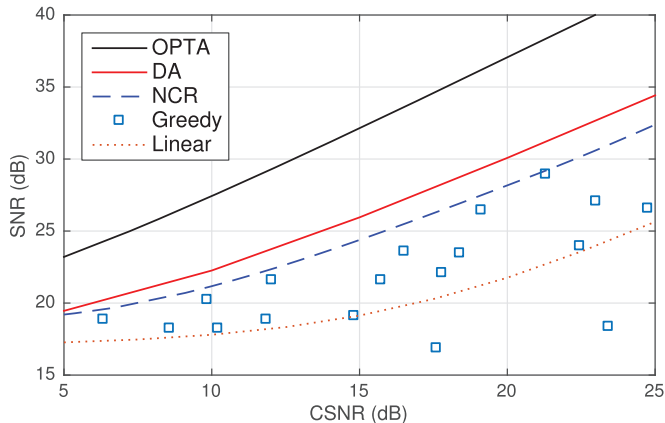


Fig. 6. The performance comparison for the side information setting, the proposed method versus the noisy relaxation (NCR), greedy optimization and the linear mappings. $\rho = 0.99$.

B. Distributed Coding Setting

In these experiments the sources are jointly Gaussian with unit variance and their correlation coefficient is denoted by ρ . We first analyze the case of individual reconstructions, where the cost is as defined in (4). The weighing coefficient η in (4) is 1.

The encoding mappings observed are many-to-one, where an example is given in Figure 7a to gain intuition into the workings of these coding schemes. From Figure 7a, where both encoders are plotted together, we see that in different source intervals, one of the mappings is many-to-one while the other one is one-to-one (usually linear). For instance, in the interval $X \in [-0.3, 0.5]$, g_1 is approximately linear while g_2 is many-to-one. Intuitively, in each of these intervals, one channel is

used as side information to reduce the uncertainty about the interval of other source.

Next, we analyze how the channel space is filled. We plot g_1 vs. g_2 in Figure 7b, which would be the channel space mapping if the two sources were fully correlated ($\rho = 1$). In case of lower correlation, a line widens into a strip (see figures and discussion in [13]), however the plot in Figure 7b is sufficient for demonstration. This mapping has the same characteristics with that of Archimedean spirals used in literature (example plotted in Figure 7c), in the sense that most likely source values are mapped to the area around origin and the mapping continues outwards in a circular fashion, to fill the channel space while preserving transmission power. In fact, spirals are suggested since they have this characteristic. Although our mappings have the same characteristic, they are far different from a spiral.

Spiral-like channel filling may sometimes be sub-optimal. The channel space can be filled in a different way, especially in case of unequal transmission powers. In Figure 8, we provide such mappings where we still see the same characteristics mentioned earlier (both sources acting as side information in different intervals), but the channel space is filled differently. Other examples can be found in literature as well, see, e.g., [12], [13].

In [13], the authors noted that for $0 < \rho < 0.95$, their structured solutions does not improve over linear solutions at high CSNR. We provide an example of non-linear scheme in Figure 9 for $\rho = 0.9$ that improves over linear solution. For lower correlations our method produces linear solutions. Based on these experiments, we reach to a similar conclusion that optimal mappings are non-linear only at high correlation - albeit our method offers non-linear gains over a larger range of ρ values.

Performance comparison of different numerical optimization techniques (DA, NCR and greedy descent with random initialization) for total power allocation case ($\lambda_1 = \lambda_2$) is provided in Figure 10a where we define $\text{SNR} = 10 \log_{10}(2/D)$ (average distortion in dB) and $\text{CSNR} = 10 \log_{10}((P(g_1) + P(g_2))/2)$ (average power in dB). We note that since individual powers are not constrained, different transmission powers are allowed in this comparison for all methods.

We also provide comparison to other coding schemes found in the literature. In [13], authors analyze parametric mappings of two types, spirals and sawtooth mappings, in distributed coding setting and compare to distributed quantizer scheme analyzed in [12]. In their comparison they use same power allocation for both encoders, as opposed to a total power allocation we consider. We therefore obtain solutions that allocate same power to both encoders. In Figure 10b, we provide comparison with our results to the ones reported in [13] for the same setting. As expected, mappings optimized in function space perform better than parametric mappings which only approximately model optimal mappings as demonstrated in Figure 7.

We finally take a look at the function computation problem for which the cost is given in (5). As a test case, we employed the difference function, $\gamma = X_1 - X_2$. Encoder mappings optimized with DA are given in Figure 11a. Both sources are mapped in many-to-one fashion with no way to resolve the uncertainty about the source interval. This is unlike previous mappings, where the uncertainty about source interval is

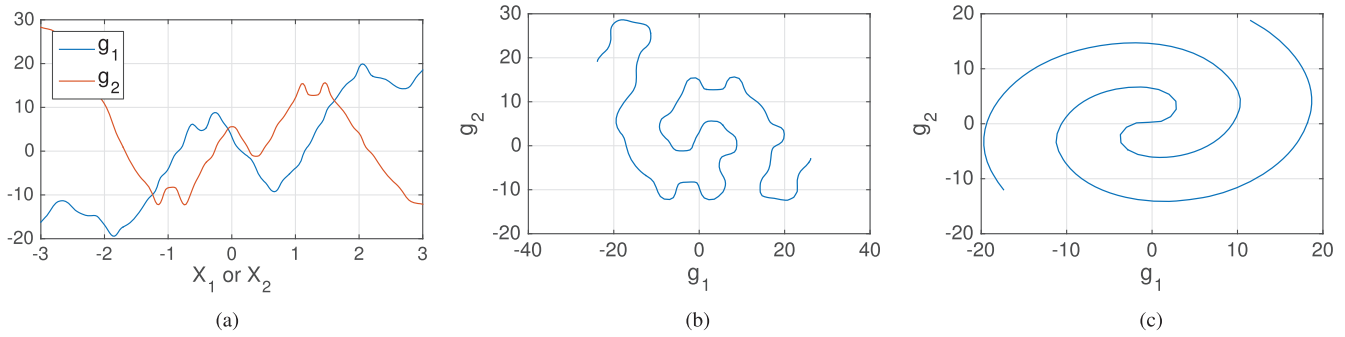


Fig. 7. Example encoding scheme for distributed coding scheme with $\rho = 0.999$. In (a), g_1 and g_2 are plotted together. In (b) we see how channel space is filled. In (c) a typical Archimedean spiral used in literature is shown.

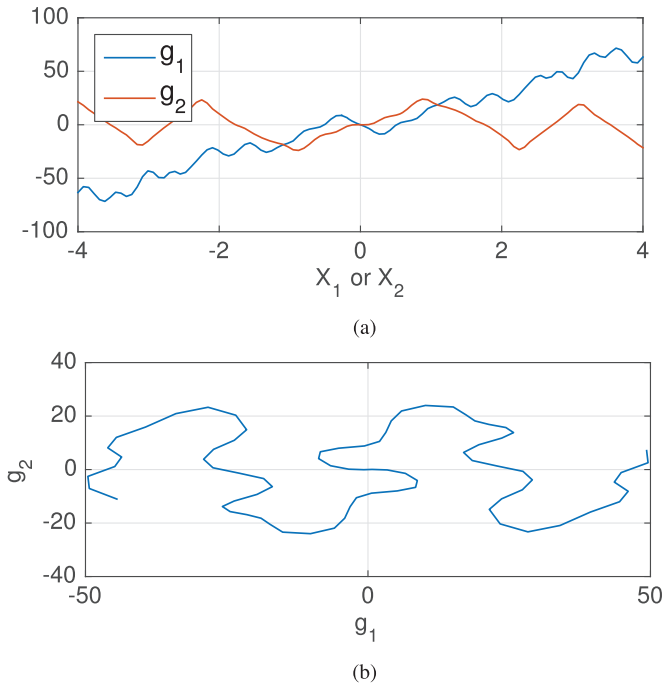


Fig. 8. An example, obtained by DA, with different transmission power constraints on encoders. (a) Both encoders are plotted together. (b) Channel space filling is shown. Although similar characteristics are observed, the channel space is filled differently.

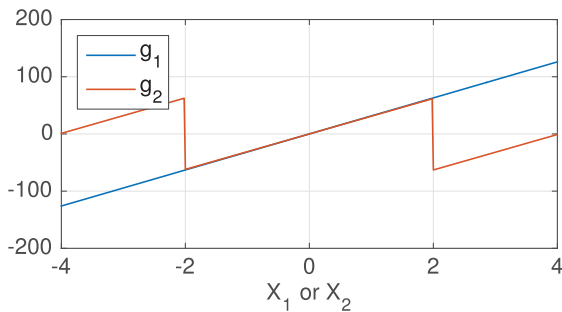


Fig. 9. Non-linear solution that improves over linear for $\rho = 0.9$. CSNR = 29 dB, SNR = 29.82 dB. Linear solution at same CSNR achieves SNR = 29.60 dB.

resolved by side information (in the distributed coding case, the other source would act as side information, at least locally). In the case of difference function, the actual values are not

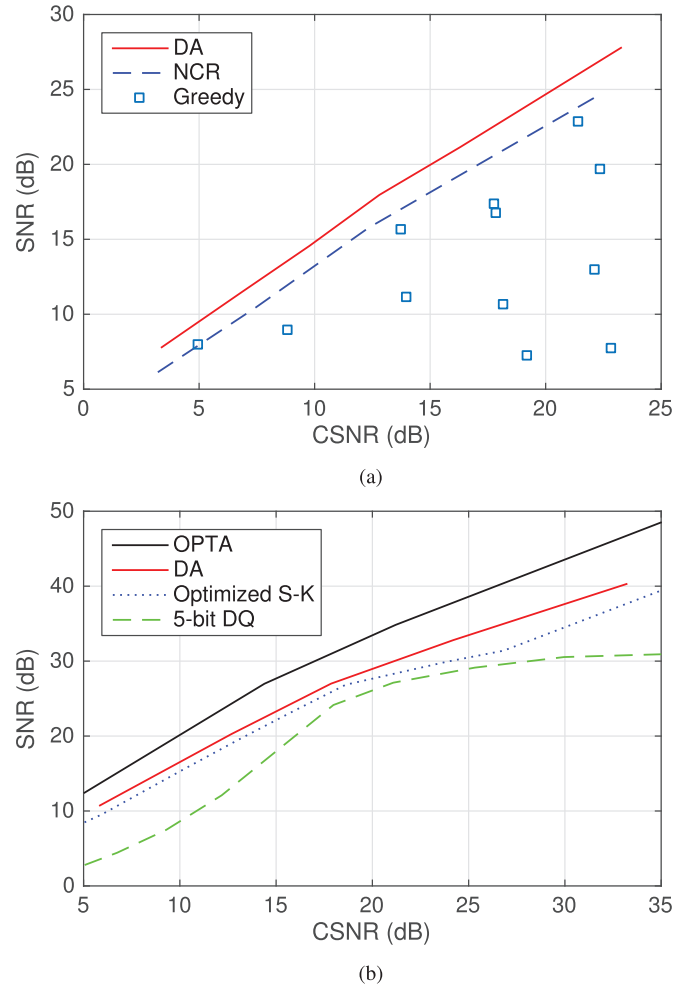


Fig. 10. (a) Performance comparison of different numerical optimization methods for distributed coding setting with the constraint on total transmission power. $\rho = 0.99$. (b) Performance comparison for distributed coding setting with other approached found in literature. Optimized S-K refers to performance of structured mappings in [13] (spirals and sawtooth mappings) and 5-bit DQ is from [12]. 5-bit DQ is optimized for 18 dB CSNR. $\rho = 0.999$.

needed, thus, both sources are mapped in many-to-one fashion. Nevertheless, the decoder is able to estimate the difference of sources accurately.

We give performance comparison in Table I where $CSNR_i = 10 \log_{10}(P(g_i))$ for $i = 1, 2$ and $SNR = 10 \log_{10}(1/D)$. DA

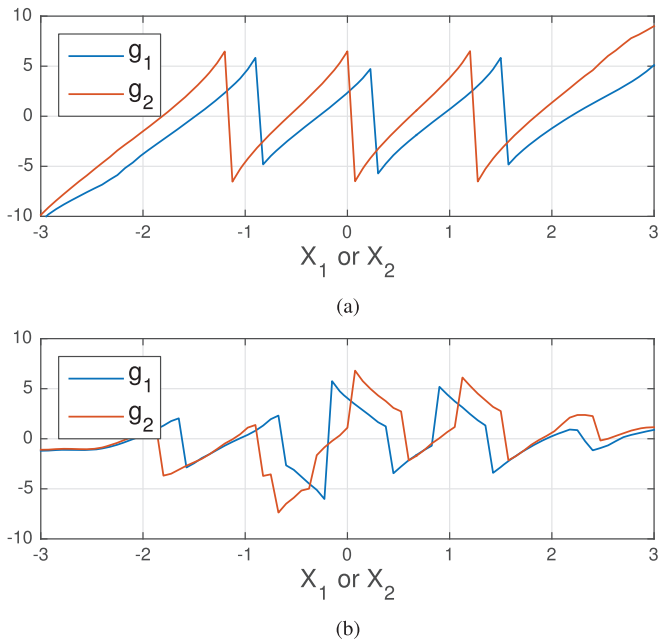


Fig. 11. Example solutions obtained for function computation problem, where $\gamma = X_1 - X_2$. (a) DA result (b) NCR result. CSNR and SNR values are in Table I.

TABLE I

PERFORMANCE OF OBTAINED MAPPINGS FOR DIFFERENCE FUNCTION

Method	CSNR ₁ (dB)	CSNR ₂ (dB)	SNR (dB)
DA	19.9	21.4	27.3
Linear-1	19.9	21.4	17.0
Linear-2	28.9	30.4	27.2
NCR	19.9	21.5	24.0

achieves 10 dB higher SNR than the linear solution with the same power allocation, whereas the linear solution that achieves the same SNR requires 9 dB more power for each channel. Although the improvements depend on the problem parameters, these results nevertheless demonstrate that there are significant gains in utilizing non-linear encoder functions instead of linear ones. DA performance is better than NCR as well, as the shape of encoders are better optimized as can be seen in comparison in Figure 11.

VI. CONCLUSIONS

In this paper, we studied the problem of finding globally optimal encoder and decoder pairs in zero delay source-channel coding, focusing on two basic network settings. Since the cost surface is riddled with local optima, we developed a method based on the deterministic annealing to approach global optimality. The numerical results show that, by using carefully optimized non-linear (and in many cases many-to-one) mappings, significant gains can be obtained over linear solutions, which are optimal in point-to-point settings (for the specific case of Gaussians under MSE). Simulation results demonstrate the performance of the proposed algorithm, which consistently outperform greedy optimization methods and noisy channel relaxation, as well as the previous approaches found in literature.

APPENDIX

PROOF OF THEOREM 1

We write the Lagrangian cost in (25) as

$$F = \sum_k \int_{\mathbf{x}} J_k(\mathbf{x}) p(k|\mathbf{x}) f_X(\mathbf{x}) d\mathbf{x} - \lambda P_E + T \sum_k \int_{\mathbf{x}} p(k|\mathbf{x}) \log p(k|\mathbf{x}) f_X(\mathbf{x}) d\mathbf{x} + TH_0, \quad (39)$$

where $J_k(\mathbf{x})$ is given in (27). From (39) it can be seen that F is convex in $p(k|\mathbf{x})$, since first term is linear and second term is convex in $p(k|\mathbf{x})$. To find the minimum, we set $\nabla_{p(k|\mathbf{x})} F = 0$:

$$J_k(\mathbf{x}) + T \log p(k|\mathbf{x}) + T = 0, \quad (40)$$

which yields

$$p(k|\mathbf{x}) = C e^{-(J_k(\mathbf{x}) - T)/T}. \quad (41)$$

The normalizing factor C is to ensure that

$$\sum_k p(k|\mathbf{x}) = 1. \quad (42)$$

Plugging (41) in (42), we have

$$C = \frac{1}{\sum_{k'} e^{-(J_{k'}(\mathbf{x}) - T)/T}}. \quad (43)$$

Plugging (43) in (41) yields (26).

REFERENCES

- [1] C. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, no. 1, pp. 379–423, 1948.
- [2] M. Gastpar, B. Rimoldi, and M. Vetterli, "To code, or not to code: Lossy source-channel communication revisited," *IEEE Trans. Inf. Theory*, vol. 49, no. 5, pp. 1147–1158, May 2003.
- [3] C. E. Shannon, "Communication in the presence of noise," *Proc. IRE*, vol. 37, no. 1, pp. 10–21, Jan. 1949.
- [4] V. A. Kotelnikov, *The Theory of Optimum Noise Immunity*. New York, NY, USA: McGraw-Hill, 1959.
- [5] A. Fuldseth and T. A. Ramstad, "Bandwidth compression for continuous amplitude channels based on vector approximation to a continuous subset of the source signal space," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Apr. 1997, vol. 4, pp. 3093–3096.
- [6] S. Y. Chung, "On the construction of some capacity approaching coding schemes," Ph.D. dissertation, Massachusetts Institute of Technology, Cambridge, MA, USA, 2000.
- [7] V. A. Vaishampayan and S. I. R. Costa, "Curves on a sphere, shift-map dynamics, and error control for continuous alphabet sources," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1658–1672, Jul. 2003.
- [8] T. A. Ramstad, "Shannon mappings for robust communication," *Teletronikk*, vol. 98, no. 1, pp. 114–128, 2002.
- [9] F. Hekland, P. A. Floor, and T. A. Ramstad, "Shannon-Kotelnikov mappings in joint source-channel coding," *IEEE Trans. Commun.*, vol. 57, no. 1, pp. 94–105, Jan. 2009.
- [10] P. Floor, A. Kim, N. Wernersson, T. Ramstad, M. Skoglund, and I. Balasingham, "Zero-delay joint source-channel coding for a bivariate Gaussian on a Gaussian MAC," *IEEE Trans. Commun.*, vol. 60, no. 10, pp. 3091–3102, Oct. 2012.
- [11] J. Kron, F. Alajaji, and M. Skoglund, "Low-delay joint source-channel mappings for the Gaussian MAC," *IEEE Commun. Lett.*, vol. 18, no. 2, pp. 249–252, Feb. 2014.

- [12] N. Wernersson, J. Karlsson, and M. Skoglund, "Distributed quantization over noisy channels," *IEEE Trans. Commun.*, vol. 57, no. 6, pp. 1693–1700, Jun. 2009.
- [13] P. A. Floor, A. N. Kim, T. A. Ramstad, and I. Balasingham, "Zero delay joint source channel coding for multivariate Gaussian sources over orthogonal Gaussian channels," *Entropy*, vol. 15, no. 6, pp. 2129–2161, 2013.
- [14] N. Wernersson and M. Skoglund, "Nonlinear coding and estimation for correlated data in wireless sensor networks," *IEEE Trans. Commun.*, vol. 57, no. 10, pp. 2932–2939, Oct. 2009.
- [15] A. Erdozain, P. M. Crespo, and B. Beferull-Lozano, "Multiple description analog joint source-channel coding to exploit the diversity in parallel channels," *IEEE Trans. Signal Process.*, vol. 60, no. 11, pp. 5880–5892, Nov. 2012.
- [16] I. Alustiza, P. M. Crespo, and B. Beferull-Lozano, "Analog multiple description joint source-channel coding based on lattice scaling," *IEEE Trans. Signal Process.*, vol. 63, no. 12, pp. 3046–3061, Mar. 2015.
- [17] Y. Hu, J. Garcia-Frias, and M. Lamarca, "Analog joint source-channel coding using non-linear curves and MMSE decoding," *IEEE Trans. Commun.*, vol. 59, no. 11, pp. 3016–3026, Nov. 2011.
- [18] N. Wernersson, M. Skoglund, and T. Ramstad, "Polynomial based analog source channel codes," *IEEE Trans. Commun.*, vol. 57, no. 9, pp. 2600–2606, Sep. 2009.
- [19] X. Chen and E. Tuncel, "Zero-delay joint source-channel coding using hybrid digital-analog schemes in the Wyner-Ziv setting," *IEEE Trans. Commun.*, vol. 62, no. 2, pp. 726–735, Feb. 2014.
- [20] P. A. Floor, T. A. Ramstad, and N. Wernersson, "Power constrained channel optimized vector quantizers used for bandwidth expansion," in *Proc. 4th Int. Symp. Wireless Commun. Syst. (ISWCS)*, Oct. 2007, pp. 667–671.
- [21] J. Karlsson and M. Skoglund, "Optimized low delay source channel relay mappings," *IEEE Trans. Commun.*, vol. 58, no. 5, pp. 1397–1404, May 2010.
- [22] E. Akyol, K. Rose, and T. Ramstad, "Optimized analog mappings for distributed source-channel coding," in *Proc. Data Compression Conf. (DCC)*, Mar. 2010, pp. 159–168.
- [23] S. Gadkari and K. Rose, "Robust vector quantizer design by noisy channel relaxation," *IEEE Trans. Commun.*, vol. 47, no. 8, pp. 1113–1116, Aug. 1999.
- [24] M. S. Mehmetoglu, E. Akyol, and K. Rose, "A deterministic annealing approach to optimization of zero-delay source-channel codes," in *Proc. IEEE Inf. Theory Workshop (ITW)*, Sep. 2013, pp. 1–5.
- [25] M. S. Mehmetoglu, E. Akyol, and K. Rose, "Optimization of zero-delay mappings for distributed coding by deterministic annealing," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, May 2014, pp. 4259–4263.
- [26] M. Goldenbaum and S. Stanczak, "Robust analog function computation via wireless multiple-access channels," *IEEE Trans. Commun.*, vol. 61, no. 9, pp. 3863–3877, Sep. 2013.
- [27] A. Giridhar and P. R. Kumar, "Toward a theory of in-network computation in wireless sensor networks," *IEEE Commun. Mag.*, vol. 44, no. 4, pp. 98–107, Apr. 2006.
- [28] A. Orlitsky and J. R. Roche, "Coding for computing," *IEEE Trans. Inf. Theory*, vol. 47, no. 3, pp. 903–917, Mar. 2001.
- [29] V. Misra, V. Goyal, and L. Varshney, "Distributed scalar quantization for computing: High-resolution analysis and extensions," *IEEE Trans. Inf. Theory*, vol. 57, no. 8, pp. 5298–5325, Aug. 2011.
- [30] B. Nazer and M. Gastpar, "Computation over multiple-access channels," *IEEE Trans. Inf. Theory*, vol. 53, no. 10, pp. 3498–3516, Oct. 2007.
- [31] K. Rose, E. Gurewitz, and G. Fox, "Statistical mechanics and phase transitions in clustering," *Phys. Rev. Lett.*, vol. 65, no. 8, pp. 945–948, 1990.
- [32] K. Rose, E. Gurewitz, and G. Fox, "Vector quantization by deterministic annealing," *IEEE Trans. Inf. Theory*, vol. 38, no. 4, pp. 1249–1257, Jul. 1992.
- [33] A. V. Rao, D. Miller, K. Rose, and A. Gersho, "A deterministic annealing approach for parsimonious design of piecewise regression models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 2, pp. 159–173, Feb. 1999.
- [34] K. Rose, "Deterministic annealing for clustering, compression, classification, regression, and related optimization problems," *Proc. IEEE*, vol. 86, no. 11, pp. 2210–2239, Nov. 1998.
- [35] E. Akyol, K. Viswanatha, K. Rose, and T. Ramstad, "On zero delay source-channel coding," *IEEE Trans. Inf. Theory*, vol. 60, no. 12, pp. 7473–7489, Dec. 2014.
- [36] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Hoboken, NJ, USA: Wiley, 1991.
- [37] A. Wyner and J. Ziv, "The rate-distortion function for source coding with side information at the decoder," *IEEE Trans. Inf. Theory*, vol. 22, no. 1, pp. 1–10, Jan. 1976.
- [38] 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.



Mustafa Said Mehmetoglu (S'12) received the B.S. degree in electrical and electronics engineering from Bilkent University, Ankara, Turkey, and the M.S. degree in electrical and computer engineering from the University of California at Santa Barbara, USA, in 2011 and 2013, respectively. He is currently pursuing the Ph.D. degree at the University of California at Santa Barbara. His research interests include signal processing, communications, optimization theory, and information theory. He is currently working on novel low-delay joint source-channel coding

approaches in communications. He was the recipient of several fellowships for his undergraduate studies due to ranking 19th in the university entrance exam nationwide, and also the recipient of the UCSB Dissertation Fellowship in 2015.



Emrah Akyol (S'03–M'12) received the Ph.D. degree from the University of California at Santa Barbara, USA, in 2011. From 2006 to 2007, he held positions at Hewlett-Packard Laboratories and NTT Docomo Laboratories, both in Palo Alto, CA, USA, where he worked on topics in video compression and streaming. From 2013 to 2014, he was a Postdoctoral Researcher with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA, USA. Currently, he is a Postdoctoral Research Associate with the

Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL, USA. His research interests include interplay of networked information theory, information economics, signal processing, communications, and stochastic control. He was the recipient of the 2010 UCSB Dissertation Fellowship, the 2014 USC Postdoctoral Training Award and was an invited participant of the 2015 NSF Early-Career Investigators Workshop on Cyber-Physical Systems and Smart City.



Kenneth Rose (S'85–M'91–SM'01–F'03) received the Ph.D. degree from the California Institute of Technology, in 1991. He then joined the Department of Electrical and Computer Engineering, University of California at Santa Barbara, USA, where he is currently a Professor. His main research activities are in the areas of information theory, communications, and signal processing, and include rate-distortion theory, source and source-channel coding, audio and video coding and networking, pattern recognition, and non-convex optimization. He is interested in the relations

between information theory, estimation theory, and statistical physics, and their potential impact on fundamental and practical problems in diverse disciplines.

Dr. Rose was the co-recipient of the 1990 William R. Bennett Prize Paper Award of the IEEE Communications Society, and the 2004 and 2007 IEEE Signal Processing Society Best Paper Awards.