

# Design of Optimal Scalar Quantizer for Sequential Coding of Correlated Sources

Huihui Wu, *Student Member, IEEE* and Sorina Dumitrescu, *Senior Member, IEEE*

**Abstract**—This work addresses the design of a sequential scalar quantizer (SSQ) for finite-alphabet correlated sources in the fixed-rate (FR) and entropy-constrained (EC) cases. The optimization problem is formulated as the minimization of a weighted sum of distortions and rates. The proposed solution is globally optimal for the class of SSQs with convex cells and is based on solving the minimum-weight path (MWP) problem in the EC case, respectively, a length-constrained MWP problem in the FR case, in a series of weighted directed acyclic graphs. The asymptotic time complexity is  $O(K_1^2 K_2^2)$ , where  $K_1$  and  $K_2$  are the respective sizes of the alphabets of the two sources. Additionally, it is proved that, by applying the proposed algorithms to discretizations of correlated sources with continuous joint probability density function, the performance approaches that of the optimal EC-SSQ, respectively FR-SSQ, with convex cells for the original sources as the accuracy of the discretization increases. Extensive experiments performed with correlated Gaussian sources validate the effectiveness in practice of the proposed approach in approximating the optimal SSQ for the case of continuous-alphabet sources.

**Index Terms**—Sequential coding, scalar quantization, globally optimal algorithm, minimum-weight path problem.

## I. INTRODUCTION

The problem of sequential coding of correlated sources (SCCS) in the information theoretical sense was introduced in [1]. The authors of [1] gave a complete characterization of the achievable rate-distortion region.

Figure 1 illustrates the framework of SCCS, where  $(X, Y)$  is a pair of jointly distributed random variables (RVs). Encoder 1 observes only the source  $X$  and encodes it at rate  $R_1$ . Decoder 1 receives the output of encoder 1 and reconstructs an estimate  $\hat{X}$  of  $X$ . Encoder 2 observes both  $X$  and  $Y$  and generates a description of  $Y$  at rate  $R_2$ . Decoder 2 utilizes the outputs of both encoders to reconstruct an estimate  $\hat{Y}$  of source  $Y$ . Note that the problem of SCCS can be regarded as a generalization of the successive refinement coding problem [3] since it reduces to the latter when the two sources coincide. In practice, the SCCS problem can be utilized to model a video sequence, where a sequence of frames corresponds to a sequence of correlated sources [1]. Moreover, it also provides a theoretical model for video compression using frame-differencing, as the encoding of a later frame refers to a previous frame [2].

In this work, we address the problem of designing a practical coding scheme for the SCCS problem which uses scalar quantization at each encoder. Specifically, encoder 1 consists of a unique scalar quantizer for the source  $X$ , while

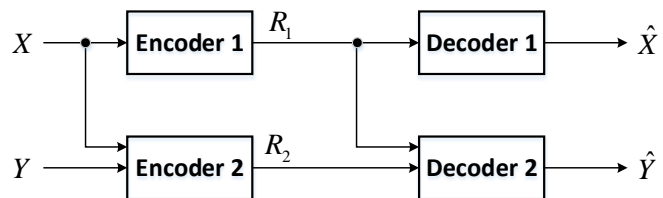


Fig. 1: Block diagram of a sequential code for correlated sources.

encoder 2 consists of a set of scalar quantizers for the source  $Y$ , each quantizer corresponding to a particular output of encoder 1. We refer to such a scheme using the term sequential scalar quantizer (SSQ). Past work on the design of SSQ includes [4] and [5], where only the fixed-rate (FR) case is considered and the quantizers are derived based on the asymptotic quantization theory. Specifically, the authors of [4] find closed form expressions for the distortion resulting from SSQ as a function of the quantizer design parameters and find the optimum parameter values that minimize the distortion. The proposed SSQ technique is utilized for color palette design of RGB images in [6], whereas an initial SSQ structure has to be preset in order to obtain the optimal number of quantization levels. It is worth pointing out that the optimization in [4] is greedy. Further, the authors of [5] improve the performance of the design procedure by considering the distribution of the unquantized scalars as well.

The most popular design approach for scalar quantizer systems is the iterative approach in the spirit of Lloyd's algorithm [7], also termed the generalized Lloyd approach. It consists of iteratively optimizing the decoder (respectively, the encoder) while the encoder (respectively, the decoder) is kept fixed. However, this design technique can only guarantee a locally optimal solution in general. Its global optimality was established so far only for certain FR quantizer systems with convex cells, for certain error functions and probability distributions [8]–[11]. Note that a quantizer cell is said to be convex if it equals the intersection of the source alphabet with a convex set. The requirement of cell convexity does not preclude optimality in the case of FR single description scalar quantizers or in the case of entropy-constrained (EC) scalar quantizers situated on the lower convex hull of the set of rate-distortion pairs, but it may in other cases [12], [13].

On the other hand, for finite-alphabet sources, globally optimal design is possible using dynamic programming, for many scalar quantizer systems. Such an approach was taken in the case of single description scalar quantizers [13]–[16], Wyner-

The authors are with the Department of Electrical and Computer Engineering, McMaster University, Hamilton, Canada (Emails: wuh58@mcmaster.ca; sorina@mail.ece.mcmaster.ca). This work was supported in part by an NSERC Discovery Grant.

Ziv scalar quantizers [13], successively refinable and multiple description scalar quantizers [13], [17]–[20], [28], joint source-channel scalar quantizers with random index assignment [11], all under the constraint of cell convexity.

This paper addresses the problem of optimal SSQ design for finite-alphabet sources in both the FR and EC cases. Note that past work [4], [5] did not consider EC-SSQ, while any optimality claims for the FR-SSQ design algorithms hold only asymptotically as the rate approaches infinity.

In the EC case, we formulate the optimization problem as the minimization of a weighted sum of the distortions at the two decoders and the rates at the two encoders. We develop a globally optimal solution algorithm with respect to the class of EC-SSQs with convex cells. The proposed algorithm relies on solving the minimum-weight path (MWP) problem in a series of appropriately constructed weighted directed acyclic graphs (WDAG). The time complexity of our solution amounts to  $O(K_1^2 K_2^2)$ , where  $K_1$  and  $K_2$  are the respective sizes of the two source alphabets.

In the FR case, we fix the number of levels of the quantizer for  $X$  and formulate the optimization problem as the minimization of a weighted sum of the distortions at the two decoders and of the rate at encoder 2. The main difference between the optimization problems in the EC and FR cases stems from the fact that in the EC case, the rate of a quantizer can be written as a sum of rates corresponding to individual quantizer cells, which is not possible in the FR case. Because of this difference, the solution to the FR problem is more involved. In particular, it needs to solve length-constrained MWP problems in a series of WDAGs, rather than unconstrained MWP problems as in the EC case. Using the straightforward solution algorithm for the length-constrained MWP problems leads to a total time complexity of  $O(K_1^2 K_2^3)$ . We further show that in some of these WDAGs, the edge weights satisfy the Monge property, fact which enables the speed up of the solution by a factor of  $K_2$ . As in the EC case, the proposed algorithm is globally optimal with respect to the class of FR-SSQs with convex cells.

As mentioned earlier, in both the EC and FR cases, we design the SSQ under the constraint of cell convexity. It is important to highlight that this constraint does not preclude the optimality of the quantizers for the source  $Y$  since the design of each such quantizer reduces to the problem of optimal scalar quantizer design for the conditional probability mass function (pmf) of  $Y$  given the particular output of the quantizer for  $X$ .

We point out that, in the case of continuous-alphabet sources, it is intuitive that approximate solutions to the EC-SSQ, respectively FR-SSQ, design problem can be obtained by applying the proposed algorithm to discretizations of the original sources. Another notable contribution of this work is a theoretical proof of the fact that the SSQ obtained in this way approaches the performance of the optimal SSQ (with convex cells) for the original sources as the discretization increases in accuracy, if the sources have a continuous joint probability density function (pdf).

It is also important to discuss the relation of our work with the algorithms proposed for the design of successively refinable scalar quantizers [13], [17], [18]. In the latter problem,

the goal is also to design a quantizer for the first encoder and conditional quantizers for the second encoder. The main difference is that for the latter problem all quantizers are designed for the same source. In our scenario, the quantizers operating at the different decoders are for distinct sources. This generalization significantly complicates the problem leading to a solution algorithm of higher computational complexity.

We would like to mention that the proposed EC-SSQ design algorithm was first presented in the conference paper [22]. In the current work, we refine the description of the algorithm and include more experimental results and discussions. Additionally, this work proposes a design algorithm for the FR case and presents an important theoretical result (Theorem 1) which does not appear in [22].

This paper is organized as follows. The next section introduces the necessary definitions and notations. Section III formulates the problem of optimal EC-SSQ design and presents the proposed solution algorithm. The problem of optimal FR-SSQ design and its solution are presented in Section IV. Section V investigates the application of the proposed designs to continuous sources. Section VI shows simulation results and, finally, Section VII concludes the paper.

## II. NOTATIONS AND PROBLEM FORMULATION

This section presents the definitions and notations used throughout this work. Let  $X$  and  $Y$  be two finite-alphabet jointly distributed RVs. Let  $P_{XY}$  denote their joint pmf. The RVs  $X$  and  $Y$  take values in the alphabets  $\mathcal{X} = \{x_1, x_2, \dots, x_{K_1}\} \subseteq \mathbb{R}$ , respectively  $\mathcal{Y} = \{y_1, y_2, \dots, y_{K_2}\} \subseteq \mathbb{R}$ , where  $K_1$  and  $K_2$  are positive integers,  $x_i < x_{i+1}$ , for  $1 \leq i \leq K_1 - 1$  and  $y_j < y_{j+1}$ , for  $1 \leq j \leq K_2 - 1$ . Let  $P_X$  and  $P_Y$  denote the marginal pmfs of  $X$  and  $Y$ , respectively.

For any positive integer  $k$  denote  $I_k \triangleq \{0, \dots, k\}$  and  $E_k \triangleq \{(u, v) \in I_k^2 | u < v\}$ . For any  $(u, v) \in E_{K_1}$  let  $C_X(u, v) \triangleq (x_u, x_v] \cap \mathcal{X} = \{x_{u+1}, \dots, x_v\}$ . For any  $(m, n) \in E_{K_2}$  denote  $C_Y(m, n) \triangleq (y_m, y_n] \cap \mathcal{Y} = \{y_{m+1}, \dots, y_n\}$ . In this work, we consider quantizers with convex cells. A subset of  $\mathcal{X}$  is said to be convex if it equals  $C_X(u, v)$  for some  $(u, v) \in E_{K_1}$ , while any convex subset of  $\mathcal{Y}$  equals  $C_Y(m, n)$  for some  $(m, n) \in E_{K_2}$ .

For any positive integer  $M$ , an ascending  $M$ -sequence for  $X$  is a sequence of integer thresholds  $\mathbf{r} \triangleq (r_0, r_1, \dots, r_M)$ , such that  $r_0 = 0 < r_1 < \dots < r_{M-1} < r_M = K_1$ . Let us denote by  $\mathcal{T}_X(M)$  the set of all such sequences. Furthermore, let  $\mathcal{T}_X \triangleq \cup_{M>0} \mathcal{T}_X(M)$ . Clearly, the encoder partition of any scalar quantizer with  $M$  convex cells for the source  $X$  can be identified with the ascending  $M$ -sequence  $\mathbf{r} \in \mathcal{T}_X(M)$ , where  $C_X(r_{i-1}, r_i]$  is the  $i$ th cell, for  $1 \leq i \leq M$ . Similarly, an ascending  $M$ -sequence for  $Y$  is any sequence of integer thresholds  $\mathbf{s} = (s_0, s_1, \dots, s_M)$  such that  $s_0 = 0 < s_1 < \dots < s_{M-1} < s_M = K_2$ . We use the notation  $\mathcal{T}_Y(M)$  for the set of all ascending  $M$ -sequences for  $Y$ , and  $\mathcal{T}_Y \triangleq \cup_{M>0} \mathcal{T}_Y(M)$ . The encoder partition of any quantizer with  $M$  convex cells for the source  $Y$  can be identified with the ascending  $M$ -sequence  $\mathbf{s} \in \mathcal{T}_Y(M)$ , where  $C_Y(s_{j-1}, s_j]$  is the  $j$ th cell, for  $1 \leq j \leq M$ . In the sequel, we use interchangeably the terms ascending sequence and quantizer (or encoder) partition.



we denote them from now on by  $D_1(\mathbf{r})$ , respectively  $D_2(\mathbf{r}, \bar{\mathbf{s}})$ . By plugging (3) in (2) we obtain

$$D_1(\mathbf{r}) = \sum_{i=1}^{M_1} \sum_{x \in C_i} (x - \hat{x}(C_i))^2 P_X(x), \quad (4)$$

$$D_2(\mathbf{r}, \bar{\mathbf{s}}) = \sum_{i=1}^{M_1} \sum_{j=1}^{M_{2,i}} \sum_{y \in C_{i,j}} (y - \hat{y}(C_{i,j}|C_i))^2 \sum_{x \in C_i} P_{XY}(x, y).$$

Let  $R_1(\mathbf{r})$  denote the rate of encoder 1 and let  $R_2(\mathbf{r}, \bar{\mathbf{s}})$  be the rate of encoder 2. The expression of the rates depends on whether the quantizers are FR or EC. Therefore, from now on we will discuss the two cases separately. In the following section we formulate the problem of optimal EC-SSQ design and propose a solution algorithm. The counterpart for the FR case is addressed in Section IV.

### III. OPTIMAL EC-SSQ DESIGN ALGORITHM

Let  $I$  and  $J$  be the random variables representing the indexes output by  $f_1$ , respectively  $f_2$ . In the EC case, the rate at encoder 1 equals the entropy of  $I$ , while the rate at encoder 2 equals the conditional entropy of  $J$  conditioned on  $I$ . Thus, we have

$$R_1(\mathbf{r}) = - \sum_{i=1}^{M_1} P(C_i) \log_2 P(C_i),$$

$$R_2(\mathbf{r}, \bar{\mathbf{s}}) = - \sum_{i=1}^{M_1} \sum_{j=1}^{M_{2,i}} P(C_i, C_{i,j}) \log_2 P(C_i, C_{i,j}) + \sum_{i=1}^{M_1} P(C_i) \log_2 P(C_i), \quad (5)$$

where  $P(C_i) \triangleq \mathbb{P}[X \in C_i]$  and  $P(C_i, C_{i,j}) \triangleq \mathbb{P}[X \in C_i, Y \in C_{i,j}]$ , for  $1 \leq i \leq M_1$  and  $1 \leq j \leq M_{2,i}$ .

Let  $\mathcal{RD}_{EC}$  denote the set of all quadruples  $(R_1(\mathbf{r}), R_2(\mathbf{r}, \bar{\mathbf{s}}), D_1(\mathbf{r}), D_2(\mathbf{r}, \bar{\mathbf{s}}))$  for all possible pairs  $(\mathbf{r}, \bar{\mathbf{s}})$ . Then any point on the lower boundary of the convex hull of  $\mathcal{RD}_{EC}$  is optimal in some sense. Any such point is the solution of the minimization of a weighted sum of the distortions and rates  $\rho_1 D_1(\mathbf{r}) + \rho_2 D_2(\mathbf{r}, \bar{\mathbf{s}}) + \lambda_1 R_1(\mathbf{r}) + \lambda_2 R_2(\mathbf{r}, \bar{\mathbf{s}})$ , for some choice of positive weights  $\rho_1, \rho_2, \lambda_1$  and  $\lambda_2$ . Note that the solution of the minimization problem remains the same if all the weights are divided by  $\rho_1 + \rho_2$ . Therefore, we formulate the optimization problem as

$$\min_{M_1, \mathbf{r} \in \mathcal{T}_X(M_1), \bar{\mathbf{s}} \in \mathcal{T}_Y^{M_1}} \mathcal{F}(\mathbf{r}, \bar{\mathbf{s}}), \quad (6)$$

where

$$\mathcal{F}(\mathbf{r}, \bar{\mathbf{s}}) \triangleq \rho D_1(\mathbf{r}) + (1 - \rho) D_2(\mathbf{r}, \bar{\mathbf{s}}) \lambda_1 R_1(\mathbf{r}) + \lambda_2 R_2(\mathbf{r}, \bar{\mathbf{s}}),$$

for some fixed  $\rho, 0 < \rho < 1, \lambda_1 > 0$  and  $\lambda_2 > 0$ . We point out that the formulation of the optimization problem as a minimization of a weighted sum of distortion(s) and rate(s) was also adopted in [13], [23], [24].

Based on relations (4)-(6) we obtain that

$$\mathcal{F}(\mathbf{r}, \bar{\mathbf{s}}) = \sum_{i=1}^{M_1} \left( \rho \sum_{x \in C_i} (x - \hat{x}(C_i))^2 P_X(x) - (\lambda_1 - \lambda_2) P(C_i) \log_2 P(C_i) + \sum_{j=1}^{M_{2,i}} \left( (1 - \rho) \sum_{y \in C_{i,j}} (y - \hat{y}(C_{i,j}|C_i))^2 \sum_{x \in C_i} P_{XY}(x, y) - \lambda_2 P(C_i, C_{i,j}) \log_2 P(C_i, C_{i,j}) \right) \right).$$

In order to simplify the expression of the cost we introduce a few more notations. For each set  $C \subseteq \mathcal{X}$  and  $C' \subseteq \mathcal{Y}$  denote

$$d_X(C) \triangleq \rho \sum_{x \in C} (x - \hat{x}(C))^2 P_X(x),$$

$$h_X(C) \triangleq -(\lambda_1 - \lambda_2) P(C) \log_2 P(C),$$

$$d_Y(C'|C) \triangleq (1 - \rho) \sum_{y \in C'} (y - \hat{y}(C'|C))^2 \sum_{x \in C} P_{XY}(x, y),$$

$$h_Y(C'|C) \triangleq -\lambda_2 P(C, C') \log_2 P(C, C').$$

Using the above notations the cost function in (6) becomes

$$\mathcal{F}(\mathbf{r}, \bar{\mathbf{s}}) = \sum_{i=1}^{M_1} \left( d_X(C_i) + h_X(C_i) + \underbrace{\sum_{j=1}^{M_{2,i}} (d_Y(C_{i,j}|C_i) + h_Y(C_{i,j}|C_i))}_{\tau(C_i, \mathbf{s}_i)} \right).$$

By examining the cost  $\mathcal{F}(\mathbf{r}, \bar{\mathbf{s}})$  we notice that for each  $i$  the contribution of the partition  $\mathbf{s}_i$  to the cost function depends on cell  $C_i$ , but does not depend on any other cell of the quantizer for  $X$ . Therefore, we will denote it by  $\tau(C_i, \mathbf{s}_i)$ . We conclude that when the partition  $\mathbf{r}$  is fixed the optimization of the partition  $\mathbf{s}_i$  can be performed separately for each  $i$ . In other words, the following holds

$$\min_{M_1, \mathbf{r} \in \mathcal{T}_X(M_1), \bar{\mathbf{s}} \in \mathcal{T}_Y^{M_1}} \mathcal{F}(\mathbf{r}, \bar{\mathbf{s}}) = \min_{M_1, \mathbf{r} \in \mathcal{T}_X(M_1)} \sum_{i=1}^{M_1} \left( d_X(C_i) + h_X(C_i) + \min_{M_{2,i}, \mathbf{s}_i \in \mathcal{T}_Y(M_{2,i})} \tau(C_i, \mathbf{s}_i) \right).$$

Further, for each  $(u, v) \in E_{K_1}$ , denote by  $\omega(C_X(u, v))$  the minimum value of  $\tau(C_i, \mathbf{s}_i)$  over all partitions  $\mathbf{s}_i$  when  $C_i = C_X(u, v]$ , in other words

$$\omega(C_X(u, v)) \triangleq \min_{M_{2,i}, \mathbf{s}_i \in \mathcal{T}_Y(M_{2,i})} \tau(C_X(u, v], \mathbf{s}_i). \quad (7)$$

With the above notation, problem (6) becomes equivalent to

$$\min_{M_1, \mathbf{r} \in \mathcal{T}(M_1)} \hat{\mathcal{F}}(\mathbf{r}) \triangleq \sum_{i=1}^{M_1} (d_X(C_i) + h_X(C_i) + \omega(C_i)). \quad (8)$$

We will show that the above problem is equivalent to an MWP problem. Indeed, consider the WDAG  $G_X(w)$ , where, for each  $(u, v) \in E_{K_1}$ ,  $w(u, v)$  is defined by

$$w(u, v) \triangleq d_X(C_X(u, v]) + h_X(C_X(u, v]) + \omega(C_X(u, v]). \quad (9)$$

Then any partition  $\mathbf{r} \in \mathcal{T}_X(M_1)$  is in a one-to-one correspondence with an  $M_1$ -edge path in  $G_X(w)$ , from the source to the final node. Additionally, the weight of the path equals the



cost  $\hat{\mathcal{F}}(\mathbf{r})$ . This implies that problem (8) is equivalent to the MWP problem in  $G_X(w)$ .

In order to solve the MWP problem in  $G_X(w)$ , we need to be able to evaluate each edge weight. Therefore, we need to solve first problem (7) for each edge  $(u, v)$ . It turns out that problem (7) is also equivalent to an MWP problem in some other WDAG. Indeed, consider the WDAG  $G_Y(w_{u,v})$ , where for each edge  $(m, n) \in E_{K_2}$ , the weight  $w_{u,v}(m, n)$  is defined as

$$w_{u,v}(m, n) \triangleq d_Y(C_Y(m, n)|C_X(u, v)) + h_Y(C_Y(m, n)|C_X(u, v)). \quad (10)$$

Then any partition  $\mathbf{s} \in \mathcal{T}_Y(M_2)$  is in a one-to-one correspondence with an  $M_2$ -edge path from the source to the final node in WDAG  $G_Y(w_{u,v})$ . The weight of the path equals the cost function in (7). Thus, problem (7) is equivalent to the MWP path problem in  $G_Y(w_{u,v})$ .

Notice that solving the MWP problem in some WDAG requires  $O(|V| + |E|)$  operations, if the weight of each edge can be evaluated in constant time, where  $V$  denotes the vertex set and  $E$  denotes the edge set. In order to enable the evaluation in constant time of each edge weight, we include a preprocessing step which computes and stores the following cumulative values

$$\begin{aligned} \varphi_{k,X}(u) &\triangleq \sum_{i=1}^u x^k P_X(x_i), \\ \varphi_{k,XY}(u, m) &\triangleq \sum_{j=1}^m \sum_{i=1}^u y^k P_{XY}(x_i, y_j), \end{aligned}$$

for  $k = 0, 1, 2$ ,  $0 \leq u \leq K_1$  and  $0 \leq m \leq K_2$ . All the above values can be computed in  $O(K_1 K_2)$  time, while the amount of memory needed store all of them is also  $O(K_1 K_2)$ . Then  $P(C_X(u, v), C_Y(m, n))$  can be computed in constant time as follows

$$P(C_X(u, v), C_Y(m, n)) = \varphi_{0,XY}(v, n) - \varphi_{0,XY}(v, m) - \varphi_{0,XY}(u, n) + \varphi_{0,XY}(u, m).$$

Similarly, the quantity  $\sum_{j=m+1}^n \sum_{i=u+1}^v y_j P_{XY}(x_i, y_j)$  can be evaluated in constant time using  $\varphi_{1,XY}(\cdot, \cdot)$ , leading further to the evaluation of  $\hat{y}(C_Y(m, n)|C_X(u, v))$  in  $O(1)$  time as well. Next notice that

$$\begin{aligned} &\sum_{j=m+1}^n (y_j - \hat{y}(C_Y(m, n)|C_X(u, v)))^2 \sum_{i=u+1}^v P_{XY}(x_i, y_j) = \\ &\sum_{j=m+1}^n \sum_{i=u+1}^v y_j^2 P_{XY}(x_i, y_j) - \\ &\hat{y}(C_Y(m, n)|C_X(u, v))^2 P(C_X(u, v), C_Y(m, n)), \end{aligned}$$

where  $\sum_{j=m+1}^n \sum_{i=u+1}^v y_j^2 P_{XY}(x_i, y_j)$  can also be computed in  $O(1)$  time based on  $\varphi_{2,XY}(\cdot, \cdot)$ .

Let us summarize now the solution algorithm to problem (6). After performing the preprocessing step the algorithm proceeds in two stages as follows.

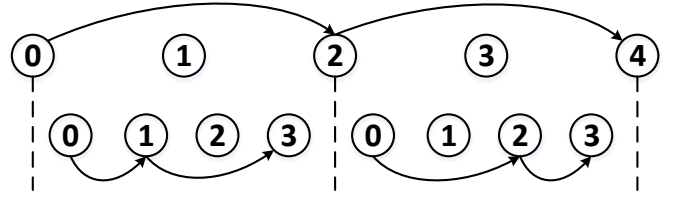


Fig. 3: The paths of graphs  $G_X(w)$  and  $G_Y(w_{u,v})$  in Example 2. Vertex es are depicted with circles and edges with arcs. The path of  $G_X(w)$  corresponds to the quantizer for  $X$  with cells  $C_1 = \{x_1, x_2\}$  and  $C_2 = \{x_3, x_4\}$ . The path in  $G_Y(w_{0,2})$  represents the quantizer for  $Y$  with cells  $C_{1,1} = \{y_1\}$  and  $C_{1,2} = \{y_2, y_3\}$ . The path in  $G_Y(w_{2,4})$  represents the quantizer for  $Y$  with cells  $C_{2,1} = \{y_1, y_2\}$  and  $C_{2,2} = \{y_3\}$ .

- 1) For each pair  $(u, v) \in E_{K_1}$ , solve the MWP problem in  $G_Y(w_{u,v})$ , where  $w_{u,v}$  is given in (10). This takes  $O(K_2^2)$  operations for each pair  $(u, v)$ . Doing so for all  $(u, v) \in E_{K_1}$  amounts to  $O(K_1^2 K_2^2)$  operations.
- 2) Solve the MWP problem in  $G_X(w)$ , where  $w$  is given in (9). This can be done in  $O(K_1^2)$  time.

In conclusion, the overall time complexity of the proposed algorithm is  $O(K_1^2 K_2^2)$ .

**Example 2.** Figure 3 shows an example of a path in  $G_X(w)$  (top) and a path in  $G_Y(w_{u,v})$  (bottom) for each edge  $(u, v)$  of the path in  $G_X(w)$ . Here  $K_1 = 4$  and  $K_2 = 3$  as in Example 1. The vertexes in each graph are represented with circles and the edges are represented with arcs. The path in  $G_X(w)$  corresponds to the quantizer for  $X$  with cells  $C_1 = \{x_1, x_2\}$  and  $C_2 = \{x_3, x_4\}$ . The path in  $G_Y(w_{0,2})$  contains edges  $(0, 1)$  and  $(1, 3)$  and corresponds to the quantizer for  $Y$  with cells  $C_{1,1} = \{y_1\}$  and  $C_{1,2} = \{y_2, y_3\}$ . The path in  $G_Y(w_{2,4})$  contains edges  $(0, 2)$  and  $(2, 3)$  and represents the quantizer for  $Y$  with cells  $C_{2,1} = \{y_1, y_2\}$  and  $C_{2,2} = \{y_3\}$ . On the other hand, we point out that the number of edges in the optimal path in  $G_Y(w_{u,v})$  may differ for different pairs  $(u, v) \in E_{K_1}$ .

#### IV. OPTIMAL FR-SSQ DESIGN ALGORITHM

In this section, we formulate the optimal FR-SSQ design problem and present its solution.

The rates in the FR case are

$$R_1(\mathbf{r}) = \log_2 M_1, \quad R_2(\mathbf{r}, \bar{\mathbf{s}}) = \sum_{i=1}^{M_1} P(C_i) \log_2 M_{2,i}. \quad (11)$$

It is easy to impose a constraint  $R_1(\mathbf{r}) \leq R_1$  on the rate of encoder 1 by fixing the number of cells in  $Q_1$  to be

$$M_1 = \lfloor 2^{R_1} \rfloor. \quad (12)$$

The problem of optimal FR-SSQ design is formulated as

$$\min_{\mathbf{r} \in \mathcal{T}_X(M_1), \bar{\mathbf{s}} \in \mathcal{T}_Y^{M_1}} \mathcal{F}'(\mathbf{r}, \bar{\mathbf{s}}), \quad (13)$$

where

$$\mathcal{F}'(\mathbf{r}, \bar{\mathbf{s}}) \triangleq \rho D_1(\mathbf{r}) + (1 - \rho) D_2(\mathbf{r}, \bar{\mathbf{s}}) + \lambda_2 R_2(\mathbf{r}, \bar{\mathbf{s}}),$$

for some fixed  $\rho, 0 < \rho < 1$ , and  $\lambda_2 > 0$ .

Let  $\mathcal{RD}_{FR}(R_1)$  denote the set of quadruples  $(R_1(\mathbf{r}), R_2(\mathbf{r}, \bar{\mathbf{s}}), D_1(\mathbf{r}), D_2(\mathbf{r}, \bar{\mathbf{s}}))$  satisfying (12). Then any point on the lower boundary of the convex hull of  $\mathcal{RD}_{FR}(R_1)$  can be obtained by solving problem (13) for some choice of  $\rho$  and  $\lambda_2$  as above.

Using the notations introduced in the previous section, the cost in (13) becomes

$$\mathcal{F}'(\mathbf{r}, \bar{\mathbf{s}}) = \underbrace{\sum_{i=1}^{M_1} \left( d_X(C_i) + \lambda_2 P(C_i) \log_2 M_{2,i} + \sum_{j=1}^{M_{2,i}} d_Y(C_{i,j}|C_i) \right)}_{\tau'(C_i, \mathbf{s}_i)}.$$

Similarly to the EC case, if cell  $C_i$  is fixed, the partition  $\mathbf{s}_i$  can be optimized by minimizing the cost  $\tau'(C_i, \mathbf{s}_i)$ . Therefore, for each  $(u, v) \in E_{K_1}$ , let us denote by  $\omega'(C_X(u, v))$  the minimum value of  $\tau'(C_i, \mathbf{s}_i)$  over all  $\mathbf{s}_i$  when  $C_i = C_X(u, v)$ , i.e.,

$$\omega'(C_X(u, v)) \triangleq \min_{M_2, \mathbf{s} \in \mathcal{J}_Y(M_2)} \tau'(C_X(u, v), \mathbf{s}). \quad (14)$$

Then problem (13) becomes equivalent to

$$\min_{\mathbf{r} \in \mathcal{J}_X(M_1)} \hat{\mathcal{F}}'(\mathbf{r}) \triangleq \sum_{i=1}^{M_1} (d_X(C_i) + \omega'(C_i)). \quad (15)$$

Consider the WDAG  $G_X(w')$ , where for each  $(u, v) \in E_{K_1}$  the weight  $w'(u, v)$  is defined as

$$w'(u, v) = d_X(C_X(u, v)) + \omega'(C_X(u, v)). \quad (16)$$

Then any ascending  $M_1$ -sequence  $\mathbf{r}$  can be identified with an  $M_1$ -edge path in  $G_X(w')$  from the source to the final node and its weight equals the cost  $\hat{\mathcal{F}}'(\mathbf{r})$ . Since the correspondence is one-to-one, it follows that problem (15) is equivalent to the  $M_1$ -edge MWP problem in  $G_X(w')$ .

In order to solve the aforementioned problem, we need to determine first the value of  $\omega'(C_X(u, v))$  by solving the minimization in (14), for each  $(u, v) \in E_{K_1}$ . Note that, unlike its counterpart (7) in the EC case, problem (14) can no longer be cast as an MWP problem. In order to solve it, notice that the following holds

$$\omega'(C_X(u, v)) = \min_{M_2} \left( \lambda_2 P(C_X(u, v)) \log_2 M_2 + \underbrace{\min_{\mathbf{s} \in \mathcal{J}_Y(M_2)} \sum_{j=1}^{M_2} d_Y(C'_j|C_X(u, v))}_{\hat{W}_{u,v}(M_2)} \right). \quad (17)$$

We conclude that the above problem can be solved in two stages.

- A) Solve first the inner minimization over ascending  $M_2$ -sequences  $\mathbf{s}$ , for each integer  $M_2 > 0$ .
- B) Solve the outer minimization over integers  $M_2 > 0$ .

For each integer  $M_2 > 0$ , the inner minimization is equivalent to the  $M_2$ -edge MWP problem in the WDAG  $G_Y(w'_{u,v})$ , where for each  $(m, n) \in E_{K_2}$ , the weight  $w'_{u,v}(m, n)$  is defined as

$$w'_{u,v}(m, n) \triangleq d_Y(C_Y(m, n)|C_X(u, v)).$$

Thus, the quantity  $\hat{W}_{u,v}(M_2)$  defined in (17) equals the weight of the  $M_2$ -edge MWP in  $G_Y(w'_{u,v})$ . As pointed out above, solving (17) can be done by determining  $\hat{W}_{u,v}(M_2)$  for each  $M_2$  and then performing a linear search over  $M_2$ .

The computation of  $\hat{W}_{u,v}(M_2)$  can be accomplished using dynamic programming (DP). The DP algorithm finds the  $k$ -edge MWP path from node 0 to node  $n$ , for each pair  $(k, n)$  with  $1 \leq k \leq M_2$  and  $1 \leq n \leq K_2$ . Let  $W_{u,v}(k, n)$  denote the weight of the  $k$ -edge MWP path from node 0 to node  $n$ . Then the following recurrence relation holds for all  $2 \leq k \leq M_2$  and  $2 \leq n \leq K_2$ ,

$$W_{u,v}(k, n) = \min_{1 \leq m < n} (W_{u,v}(k-1, m) + w'_{u,v}(m, n)). \quad (18)$$

Clearly,  $W_{u,v}(1, m) = w'_{u,v}(0, m)$  for all  $m \in I_{K_2} \setminus \{0\}$ . The DP process solves (18) for all pairs  $(k, n)$ ,  $1 \leq k \leq M_2$ ,  $1 \leq n \leq K_2$ , in lexicographical order. The value  $\hat{W}_{u,v}(M_2)$  sought of equals  $W_{u,v}(M_2, K_2)$ . The total amount of operations reaches  $O(M_2 K_2^2)$ .

Note that the above procedure to solve the  $M_2$ -edge MWP problem, also solves the  $k$ -edge MWP problem for all smaller path lengths  $k$ , for  $1 \leq k < M_2$ . Since the maximum possible value of  $M_2$  is  $K_2$ , it follows that solving the  $M_2$ -edge MWP problem for all  $1 \leq M_2 \leq K_2$  can be done in  $O(K_2^3)$  time. Since the additional linear search over  $M_2$  in (17) takes only  $O(K_2)$  time, it follows that problem (17) can be solved in  $O(K_2^3)$  time.

Next we will show that the edge weights in the WDAG  $G_Y(w'_{u,v})$  satisfy the so-called Monge property, fact which allows for a speed-up of the DP algorithm.

**Lemma.** The edge weights in the WDAG  $G_Y(w'_{u,v})$  satisfy the Monge property, i.e., the following holds

$$w'_{u,v}(m, n) + w'_{u,v}(m', n') \leq w'_{u,v}(m, n') + w'_{u,v}(m', n),$$

for all  $0 \leq m < m' < n < n' \leq K_2$ . (19)

*Proof:* Let  $C = C_X(u, v)$ ,  $P_C(y) \triangleq \frac{\sum_{x \in C} P_{XY}(x, y)}{P(C)}$  and  $\eta(m, n) \triangleq \sum_{j=m+1}^n (y_j - \hat{y}(C_Y(m, n)|C))^2 P_C(y_j)$ . Then we have

$$w'_{u,v}(m, n) = (1 - \rho) P(C_X(u, v)) \eta(m, n). \quad (20)$$

Note that  $P_C(y)$  is a pmf and  $\hat{y}(C_Y(m, n)|C) = \frac{\sum_{j=m+1}^n y_j P_C(y_j)}{\sum_{j=m+1}^n P_C(y_j)}$ . Then according to [15], [16], the function  $\eta(m, n)$  satisfies the Monge property, i.e., the following holds

$$\eta(m, n) + \eta(m', n') \leq \eta(m, n') + \eta(m', n),$$

for all  $0 \leq m < m' < n < n' \leq K_2$ . The above property in conjunction with (20) implies (19), thus completing the proof. ■

Since the weights  $w'_{u,v}(m, n)$  of the WDAG  $G_Y(w'_{u,v})$  satisfy the Monge property, the DP algorithm used to solve the problem at stage A can be sped up by a factor of  $K_2$  [15], [16]. Specifically, this is done by applying the so-called SMAWK algorithm introduced in [25] to compute all values  $W_{u,v}(k, n)$  for all  $n$  and fixed  $k$ , in  $O(K_2)$  operations. This implies that problem (17) can be solved in  $O(K_2^2)$  time. It follows that computing  $\omega'(C_X(u, v))$  for all pairs  $(u, v) \in E_{K_1}$  takes  $O(K_1^2 K_2^2)$  operations.

Let us summarize now the proposed solution to the optimal FR-SSQ design problem (13). We start with a preprocessing step as in the EC case. After that the algorithm proceeds as follows.

- 1) For each pair  $(u, v) \in E_{K_1}$ , solve problem (17) in the following two stages.
  - A) Solve the  $M_2$ -edge MWP problem in  $G_Y(w'_{u,v})$  for all  $1 \leq M_2 \leq K_2$ . To this end, for each  $1 \leq k \leq K_2$ , use SMAWK to compute  $W_{u,v}(n)$  for all  $1 \leq n \leq K_2$ .
  - B) Compute

$$w'(C_X(u, v)) = \min_{M_2} (\lambda_2 P(C_X(u, v)) \log_2 M_2 + \hat{W}_{u,v}(M_2)).$$

- 2) Solve the  $M_1$ -edge MWP problem in the WDAG  $G_X(w')$ .

Recall that the preprocessing step needs  $O(K_1 K_2)$  time. Further, Step 1 requires  $O(K_1^2 K_2^2)$  operations. Step 2 can be accomplished in  $O(M_1 K_1^2)$  running time. In conclusion, the overall running time to solve problem (13) is  $O(K_1^2 K_2^2)$  assuming that  $M_1 = O(K_2^2)$ .

## V. APPLICATION TO CONTINUOUS SOURCES

In this section we assume that the sources  $X$  and  $Y$  are continuous and apply the proposed algorithms to discretized versions of  $X$  and  $Y$ . We show that the EC-SSQ, respectively FR-SSQ, obtained in this way approaches in performance the optimal EC-SSQ, respectively FR-SSQ, with convex cells for the original sources as the discretization increases in accuracy.

First we need to introduce some notations. For any pair of real-valued RVs  $(X, Y)$  with joint pdf  $f_{XY}$ , for each positive real value  $B$  and positive integer  $K$ , we define the pair of continuous RVs  $(X_B, Y_B)$  and the pair of discrete RVs  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$  as follows.  $(X_B, Y_B)$  is the truncation of  $(X, Y)$  to the set  $[-B, B] \times [-B, B]$ , i.e., its pdf is  $f_{X_B Y_B}(x, y) \triangleq \frac{f_{XY}(x, y)}{\int_{-B}^B \int_{-B}^B f_{XY}(x, y) dx dy}$  when  $(x, y) \in [-B, B] \times [-B, B]$  and 0 otherwise. The marginal pdfs of  $X_B$  and  $Y_B$  are denoted by  $f_{X_B}$  and  $f_{Y_B}$ , respectively. Further,  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$  is the quantized version of  $(X_B, Y_B)$  using a product scalar quantizer. More specifically, each scalar quantizer has  $K$  cells of equal size, and the centroid of each cell as the reconstruction value<sup>1</sup>. Thus, the thresholds of each scalar quantizer are  $t_0^{(B)}, \dots, t_K^{(B)}$ , where  $t_k^{(B)} \triangleq -B + \frac{2kB}{K}$ ,  $0 \leq k \leq K$ . Let  $\mathcal{U}_{B,K}$  denote the set of these thresholds. The alphabet of  $\tilde{X}_{B,K}$  is  $\tilde{\mathcal{X}}_{B,K} = \{x_k^{(B)} | 1 \leq k \leq K\}$ , where  $x_k^{(B)} \triangleq \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} x f_{X_B}(x) dx / \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} f_{X_B}(x) dx$  if  $\int_{t_{k-1}^{(B)}}^{t_k^{(B)}} f_{X_B}(x) dx > 0$  and  $x_k^{(B)} \triangleq (t_{k-1}^{(B)} + t_k^{(B)})/2$  otherwise. The alphabet of  $\tilde{Y}_{B,K}$  is  $\tilde{\mathcal{Y}}_{B,K} = \{y_k^{(B)} | 1 \leq k \leq K\}$ , where  $y_k^{(B)} \triangleq \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} y f_{Y_B}(y) dy / \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} f_{Y_B}(y) dy$  if  $\int_{t_{k-1}^{(B)}}^{t_k^{(B)}} f_{Y_B}(y) dy > 0$  and  $y_k^{(B)} \triangleq (t_{k-1}^{(B)} + t_k^{(B)})/2$  otherwise. The joint pmf of  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$  is  $P_{\tilde{X}_{B,K} \tilde{Y}_{B,K}}(x_k^{(B)}, y_l^{(B)}) \triangleq \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} f_{X_B Y_B}(x, y) dy dx$ ,  $1 \leq k, l \leq K$ .

<sup>1</sup>We make the convention that, if the probability of a cell is 0, then its centroid is the middle of the corresponding interval.

An SSQ for a continuous source is specified by the encoding functions  $f_1, f_2$  and the decoding functions  $g_1, g_2$ , as in (1), where  $\mathcal{J}_1 = \{1, 2, \dots, M_1\}$  or  $\mathcal{J}_1 = \mathbb{Z}$  and  $\mathcal{J}_2 = \{1, 2, \dots, M_2\}$  or  $\mathcal{J}_2 = \mathbb{Z}$ . Note that we also consider the possibility that  $\mathcal{J}_1 = \mathbb{Z}$  and  $\mathcal{J}_2 = \mathbb{Z}$  in the EC case. The following restrictions are imposed on the partitions generated by  $f_1$  and  $f_2(i, \cdot)$ ,  $i \in \mathcal{J}_1$ .

- C1) Each partition has convex cells, thus the cells are intervals, open at the left end and closed at the right end (except when the right end equals infinity).
- C2) Each partition has a finite number of cells in any bounded interval<sup>2</sup>.

For simplicity, let us denote  $\mathbf{Q} = (f_1, f_2, g_1, g_2)$ . When applying the SSQ  $\mathbf{Q}$  to a pair of RVs  $(X', Y')$ , we denote by  $D_1(\mathbf{Q}, X')$  and  $D_2(\mathbf{Q}, X', Y')$  the distortions at the first and second decoder, respectively, i.e.,

$$D_1(\mathbf{Q}, X') \triangleq \mathbb{E}[(X' - \hat{X}')^2],$$

$$D_2(\mathbf{Q}, X', Y') \triangleq \mathbb{E}[(Y' - \hat{Y}')^2],$$

where  $\hat{Y}' = g_2(f_1(X'), f_2(f_1(X'), Y'))$  and  $\hat{X}' = g_1(f_1(X'))$ . The rates of the two encoders in the EC case will be denoted by  $R_{EC,1}(\mathbf{Q}, X')$  and  $R_{EC,2}(\mathbf{Q}, X', Y')$ , respectively. Thus,

$$R_{EC,1}(\mathbf{Q}, X') \triangleq -\mathbb{E}[\log_2 P(f_1(X'))],$$

$$R_{EC,2}(\mathbf{Q}, X', Y') \triangleq -\mathbb{E}[\log_2 P(f_2(f_1(X'), Y') | f_1(X'))],$$

where, for a discrete RV  $\tilde{Z}$ ,  $P(\tilde{Z})$  denotes its pmf, i.e.,  $P(\tilde{Z}) = P_{\tilde{Z}}(\tilde{Z})$ . The rates in the FR case are  $R_{FR,1}(\mathbf{Q}, X') \triangleq \log_2 M_1$  and  $R_{FR,2}(\mathbf{Q}, X', Y') \triangleq -\mathbb{E}[\log_2 M_{2,I}]$ . Note that in the FR case, we necessarily have  $\mathcal{J}_1$  and  $\mathcal{J}_2$  finite. We denote by  $\mathcal{Q}_{EC}$  and by  $\mathcal{Q}_{FR}(M_1)$  the class of EC-SSQs and of FR-SSQs defined as above (and, thus, satisfying conditions C1 and C2), respectively. Finally, consider fixed  $0 < \rho < 1$ ,  $\lambda_1 > 0$  and  $\lambda_2 > 0$  and denote

$$\mathcal{F}_{EC}(\mathbf{Q}, X', Y') \triangleq \rho D_1(\mathbf{Q}, X') + (1 - \rho) D_2(\mathbf{Q}, X', Y') + \lambda_1 R_{EC,1}(\mathbf{Q}, X') + \lambda_2 R_{EC,2}(\mathbf{Q}, X', Y'),$$

$$\mathcal{F}_{FR}(\mathbf{Q}, X', Y') \triangleq \rho D_1(\mathbf{Q}, X') + (1 - \rho) D_2(\mathbf{Q}, X', Y') + \lambda_2 R_{FR,2}(\mathbf{Q}, X', Y').$$

The proof of the following result is deferred to the appendix.

**Theorem 1:** Let  $(X, Y)$  be a pair of jointly distributed real-valued RVs with a continuous joint pdf  $f_{XY}$  with finite variance. For each positive real value  $B$  and positive integer  $K$ , let  $\hat{\mathbf{Q}}_{B,K}$  denote the optimal EC-SSQ with convex cells for the pair of discrete RVs  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$ . Then the following holds

$$\lim_{B \rightarrow \infty} \lim_{K \rightarrow \infty} \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) = \inf_{\mathbf{Q} \in \mathcal{Q}_{EC}} \mathcal{F}_{EC}(\mathbf{Q}, X, Y). \quad (21)$$

<sup>2</sup>Note that, in the EC case, considering only partitions where the number of cells is finite in any bounded interval does not preclude the optimality of the quantizer for  $Y$ , according to [28]. There exists the possibility that the arguments of [28] could be extended to prove a similar claim for the quantizer for  $X$ , too. The investigation of such a possibility is left for future work.

Furthermore, for each positive real value  $B$  and positive integers  $K$  and  $M_1$ , let  $\hat{\mathbf{Q}}_{B,K}(M_1)$  denote the optimal FR-SSQ with convex cells and with  $M_1$  cells in the encoder 1 partition, for the pair of discrete RVs  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$ . Then the following holds

$$\lim_{B \rightarrow \infty} \lim_{K \rightarrow \infty} \mathcal{F}_{FR}(\hat{\mathbf{Q}}_{B,K}(M_1), \tilde{X}_{B,K}, \tilde{Y}_{B,K}) = \inf_{\mathbf{Q} \in \mathcal{Q}_{FR}(M_1)} \mathcal{F}_{FR}(\mathbf{Q}, X, Y), \quad (22)$$

for each positive integer  $M_1$ .

## VI. EXPERIMENTAL RESULTS AND DISCUSSION

This section assesses the practical performance of the proposed EC-SSQ and FR-SSQ design algorithms for discretized Gaussian sources. We start with a pair  $(X, Y)$  of correlated Gaussian sources, both with 0 mean and variance 1, with joint pdf

$$f_{XY}(x, y) = \frac{1}{2\pi\sqrt{1-c^2}} \exp\left(-\frac{x^2 + y^2 - 2xyc}{2(1-c^2)}\right),$$

where  $c$  is the correlation coefficient. We consider  $c = 0.5$  and  $c = 0.9$  in this section.

Next we consider the pair of discrete sources  $(\tilde{X}, \tilde{Y}) = (\tilde{X}_{B_1, K_1}, \tilde{Y}_{B_2, K_2})$ , where  $B_1 = 3$ ,  $B_2 = 5$ ,  $K_1 = 100$  and  $K_2 = 160$ . The proposed EC-SSQ and FR-SSQ design algorithms are applied to the pair of discrete RVs  $(\tilde{X}, \tilde{Y})$  and the obtained SSQs are extended to SSQs for the continuous sources  $(X, Y)$ . Then the distortions at the two decoders, denoted by  $D_1$ , respectively  $D_2$ , and the rates of the two encoders, denoted by  $R_1$ , respectively  $R_2$ , are evaluated for the extended SSQs applied to  $(X, Y)$ .

An SSQ for the discrete sources  $(\tilde{X}, \tilde{Y})$  is extended to an SSQ for  $(X, Y)$  by extending each partition of the alphabet of  $\tilde{X}$  and each partition of the alphabet of  $\tilde{Y}$  to a partition of  $\mathbb{R}$  with the same number of cells as follows. A partition for  $\tilde{X}$  specified by the sequence of thresholds  $0 = r_0 < r_1 < \dots < r_{M_1} = K_1$  is extended to the partition of  $\mathbb{R}$  with thresholds  $(-\infty, t_{r_1}^{(B_1)}, \dots, t_{r_{M_1-1}}^{(B_1)}, \infty)$ . Likewise, a partition for  $\tilde{Y}$  specified by the sequence of thresholds  $0 = s_0 < s_1 < \dots < s_{M_2} = K_2$  is extended to the partition of  $\mathbb{R}$  with thresholds  $(-\infty, t_{s_1}^{(B_2)}, \dots, t_{s_{M_2-1}}^{(B_2)}, \infty)$ .

We first consider the case of EC-SSQ. We ran the proposed algorithm for optimal EC-SSQ design for four values of  $\rho$ , namely  $\rho = 0.1, 0.5, 0.9, 0.95$ , and for a large set of values of  $\lambda_1$  and  $\lambda_2$  with  $\lambda_1 \in [0.01, 1.50]$  and  $\lambda_2 \in [0.01, 1.0]$ .

Figures 4 and 5 illustrate the performance comparison against the theoretical rate-distortion bounds. Figures 4a and 5a plot the distortion pairs  $(D_1, D_2)$  obtained in our experiments for  $c = 0.9$  and  $c = 0.5$ , respectively. Each figure also shows the boundary of the theoretical region of nontrivial distortion pairs, which is characterized by  $0 \leq D_1 \leq 1$  and  $D_2 \leq 1 - c^2(1 - D_1)$ . These figures show that, by varying the parameters  $\rho$ ,  $\lambda_1$  and  $\lambda_2$ , the proposed design is able to achieve a dense set of distortion pairs covering fairly well the theoretical distortion region. For each distortion pair  $(D_1, D_2)$  achieved by our scheme we compute the rate-gap pair  $(\Delta R_1, \Delta R_2)$  relative to the theoretical lower bound,

namely  $\Delta R_i = R_i - R_i^*$ ,  $i = 1, 2$ , where  $(R_1^*, R_2^*)$  denotes the pair of information theoretical lower bounds on the rates at the two encoders for the distortion pair  $(D_1, D_2)$ . According to [1], we have

$$R_1^* = \frac{1}{2} \log_2 \frac{1}{D_1}, \quad R_2^* = \frac{1}{2} \log_2 \frac{1 - c^2(1 - D_1)}{D_2}.$$

The rate-gap pairs are plotted in Figures 4b and 5b for  $c = 0.9$  and  $c = 0.5$ , respectively. Note that the existence of a gap is expected since the theoretical bound is achieved using vector quantization with dimension approaching  $\infty$ , while we use scalar quantization. The rate gap between the optimum EC scalar quantizer and the rate-distortion limit was proved in [26] to be  $\frac{1}{2} \log_2 \frac{2\pi e}{12} = 0.2546$  bits/sample at high resolution.

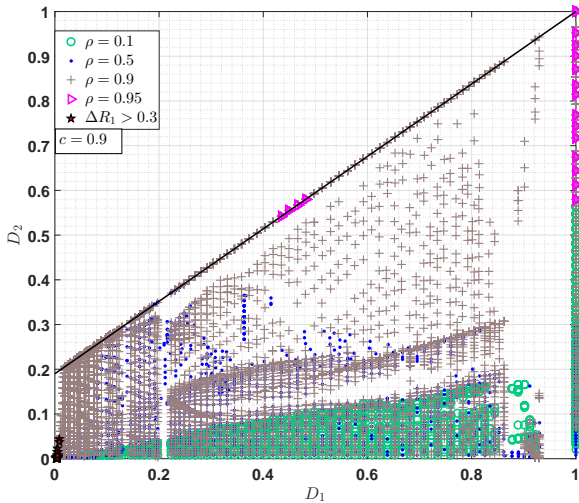
As it can be seen from Figures 4b and 5b, in most of the cases the rate-gap at encoder 2 is within 0.254 bits/sample, while the gap at encoder 1 is within 0.274 bits/sample, which is very close to the gap due to the low dimensionality of the EC-SSQ. This fact demonstrates the effectiveness of the proposed EC-SSQ design algorithm as an approximation of the optimal EC-SSQ for continuous sources.

We also mention that the largest value of  $\Delta R_2$  is only slightly higher than the benchmark value of 0.2546, namely it is 0.257 bits/sample for  $c = 0.9$ , respectively 0.262 bits/sample for  $c = 0.5$ . On the other hand, there are several cases for which the rate-gap  $\Delta R_1$  ranges between 0.3 and 0.4. The corresponding rate-gap pairs and distortion pairs are marked using star-shaped markers in Figures 4 and 5. We observe that these cases with excess rate loss are obtained when  $D_1$  is very small (thus,  $R_1$  is very high), while  $D_2 \lesssim 0.1$ . One possible reason for this additional rate loss could be the coarseness of discretization for the source  $X$ . Another possible reason could be the additional tension in the optimization of encoder 1 generated by the competing requirements at the two decoders. Namely, there is tension between ensuring a good reconstruction of the source  $X$  as well as facilitating an efficient encoder for the source  $Y$ .

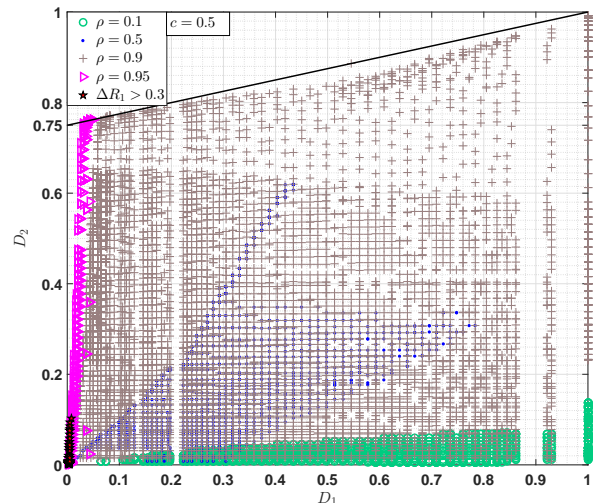
It is also interesting to investigate the impact that the refinement of the discretization has on the EC-SSQ performance. Table I compares the rate-distortion performance for four pairs  $K_1, K_2$  representing a gradual increase in the discretization accuracy. The EC-SSQ design algorithm is applied in all four cases to the same parameters, namely  $c = 0.9$ ,  $\lambda_1 = 0.22$ ,  $\lambda_2 = 0.15$  and  $\rho = 0.5$ . The pair  $K_1 = 100, K_2 = 160$  represents the coarsest discretization. The discretization is refined gradually by multiplying the initial values of  $K_1$  and  $K_2$  by two, five and ten, respectively. It can be noted that the rate gaps generally decrease, as expected, but the decrease is very small. In particular, the relative decrease of  $\Delta R_1$  from the initial to the final value is of 0.5%, while for  $\Delta R_2$  the relative decrease is of 0.28%.

It is instructive to analyze the structure of the encoder partitions generated by the proposed approach. Note that in the sequel, the distortion is represented in dB, i.e., as  $10 \log_{10} D$ . Figure 6 illustrates the optimized encoder partitions of the proposed EC-SSQ with  $R_1 = 1.3173$  and  $R_2 = 1.0430$ , when  $c = 0.9$  and  $\rho = 0.5$ . In this example, the source  $X$  is quantized to  $M_1 = 3$  cells with the sequence of thresholds

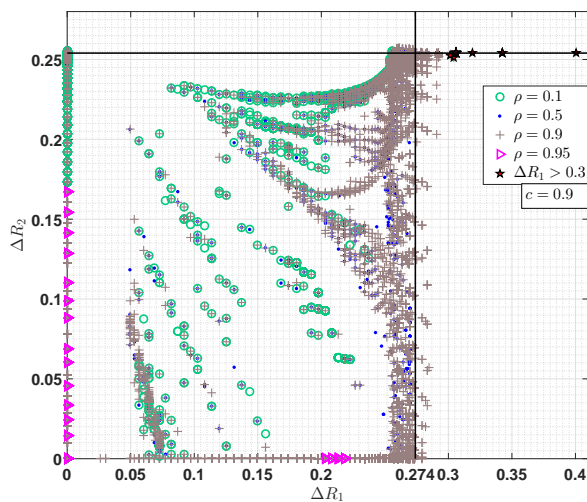




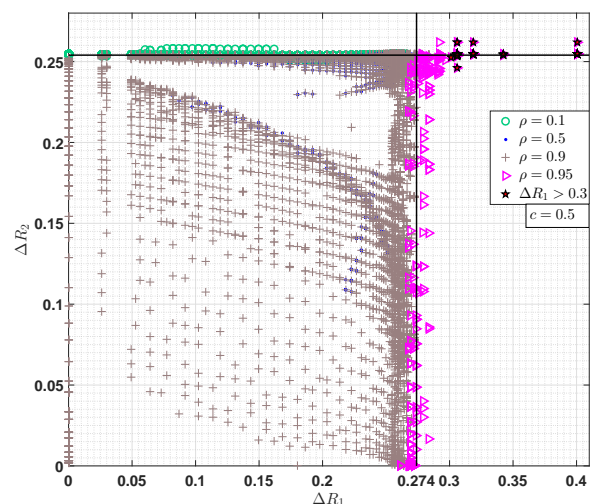
(a) Achievable distortion pairs against the theoretical bound.



(a) Achievable distortion pairs against the theoretical bound.



(b) Gap to the theoretical minimum rate.



(b) Gap to the theoretical minimum rate.

Fig. 4: Comparison between the achievable rate-distortion performance and the theoretical bound for  $c = 0.9$ .

Fig. 5: Comparison between the achievable rate-distortion performance and the theoretical bound for  $c = 0.5$ .

TABLE I: Rate-distortion performance comparison of the proposed EC-SSQ for various  $K_1$  and  $K_2$ .

$(K_1, K_2)$	(100, 160)	(200, 320)	(500, 800)	(1000, 1600)
$R_1$	1.3173	1.3030	1.3059	1.3030
$D_1$	0.2307	0.2349	0.2340	0.2349
$R_2$	1.0430	1.0461	1.0442	1.0452
$D_2$	0.1196	0.1201	0.1201	0.1202
$\Delta R_1$	0.2593	0.2580	0.2582	0.2580
$\Delta R_2$	0.2150	0.2145	0.2144	0.2144

$(-\infty, -0.9, 0.9, \infty)$ . For each  $i = 1, 2, 3$ , all quantizers of source  $Y$  have  $M_{2,i} = 4$  cells. The partitions corresponding to the quantizers for  $Y$ , for  $i = 1, 2, 3$ , are defined by the sequences of thresholds  $(-\infty, -3.125, -1.75, -0.188, \infty)$ ,  $(-\infty, -1.438, 0, 1.438, \infty)$  and  $(-\infty, 0.188, 1.75, 3.125, \infty)$ , respectively. In addition, the contour of the joint pdf  $f_{XY}$  is also plotted in Figure 6, where the probability decreases as

the color changes from green to blue. It is worth pointing out that the output of the quantizer of  $Y$  is more densely spaced where the joint probability takes on large values, as expected.

Figure 7 plots the distortion  $D_2$  of the proposed EC-SSQ, versus the rate  $R_2$ , when the pair  $(R_1, D_1)$  is fixed, for three cases of  $(R_1, D_1)$  with  $c = 0.9$  and  $\rho = 0.5$ . As expected, for a fixed pair  $(R_1, D_1)$ , the distortion at the second decoder decreases steadily as the rate at encoder 2 increases. On the other hand, when the rate  $R_2$  is kept fixed, the performance at decoder 2 also improves consistently with the increase of the rate at encoder 1. In particular, when  $R_1$  increases from 1.1983 to 1.6095, the performance at decoder 2 jumps up by about 0.9 dB. The further increase of  $R_1$  to 1.9367 leads to another gain of about 0.5 dB at decoder 2. This is expected since, intuitively, increasing  $R_1$  corresponds to refining the information about the source  $X$ . Since  $X$  and  $Y$  are correlated,

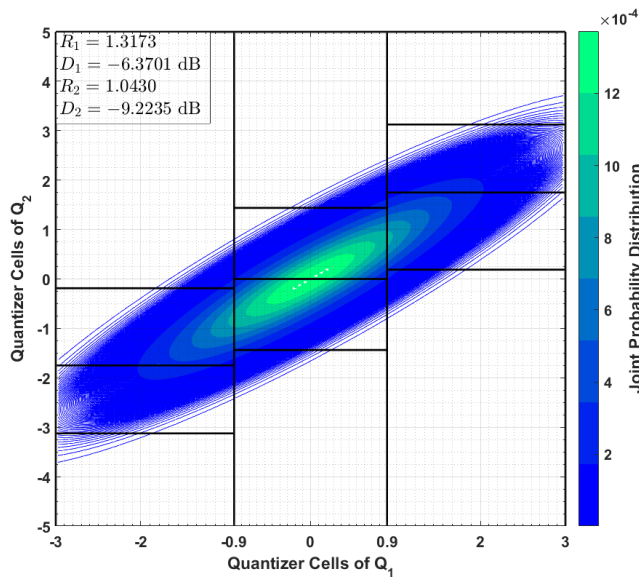


Fig. 6: Example of optimized encoder partitions of the proposed EC-SSQ, when  $c = 0.9$  and  $\rho = 0.5$ .

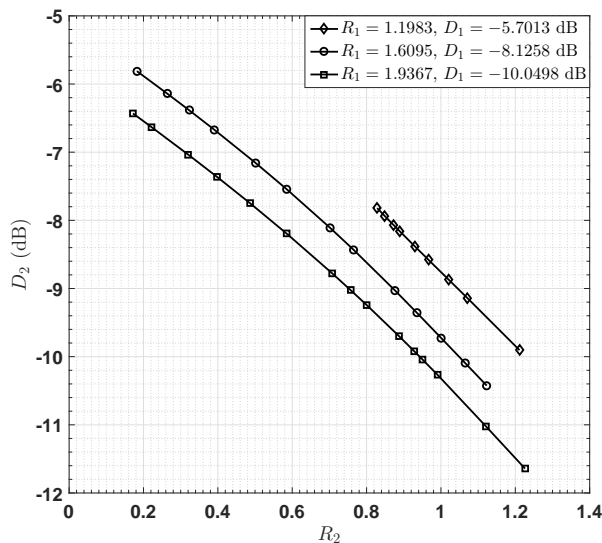


Fig. 7: Performance of proposed EC-SSQ at decoder 2, for three pairs  $(R_1, D_1)$  when  $c = 0.9$  and  $\rho = 0.5$ .

the refinement of the information about  $X$  leads to more information about  $Y$ . Thus, the rate at encoder 2 is used to further refine the information which is already available about  $Y$  through the reconstruction of  $X$ .

Next we assess the performance of the proposed FR-SSQ design algorithm in comparison with the level-constrained practical SSQ scheme developed in [4], based on the asymptotic quantization theory. The authors of [4] use the following quantizer density functions for  $Q_1$ , respectively  $Q_{2,i}$ ,

$$\lambda(x) = \frac{f_X(x)^{1/3}}{\int f_X(x)^{1/3} dx}, \quad \lambda(y|C_i) = \frac{f(y|C_i)^{1/3}}{\int f(y|C_i)^{1/3} dy},$$

to derive the asymptotical expressions of the distortions as

the rates approach infinity. Further, based on the asymptotical analysis, they propose a practical scheme operating at finite rates. Note that the design of [4] is performed under the constraint that  $\sum_{i=1}^{M_1} M_{2,i} = N$ , for some target value  $N$ . The practical construction of [4] proceeds as follows. First, the encoding function  $f_1$  partitions the real line into  $M_1$  cells such that the area under the function  $\lambda(x)$  within each cell equals  $1/M_1$ , using the marginal pdf  $f_X(x)$ . Subsequently, the values of  $M_{2,i}$  are computed using

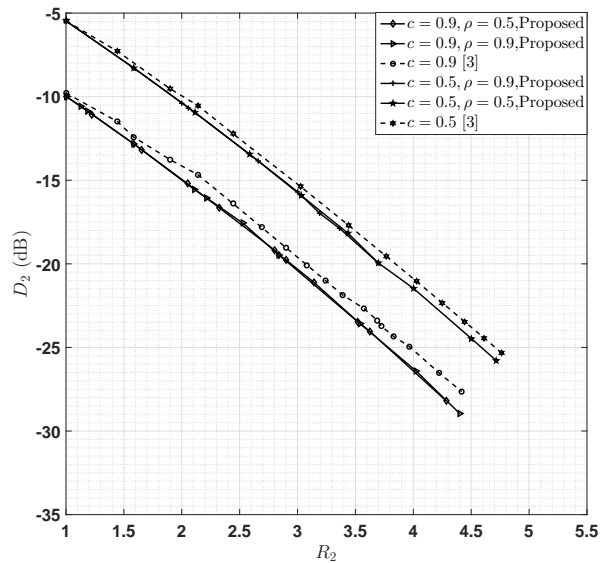
$$M_{2,i} = \left\lceil N \frac{[\|f(y|C_i)\|_{1/3} P(C_i)]^{1/3}}{\sum_{i=1}^{M_1} [\|f(y|C_i)\|_{1/3} P(C_i)]^{1/3}} \right\rceil,$$

where  $\|f(x)\|_m = [\int f(x)^m dx]^{1/m}$ , while  $\lceil \cdot \rceil$  denotes rounding to the nearest integer. Subsequently, for each cell  $i$ ,  $1 \leq i \leq M_1$ , the encoding function  $f_2(i, \cdot)$  partitions the real line into  $M_{2,i}$  cells such that the area under  $\lambda(y|C_i)$  within each cell equals  $1/M_{2,i}$ , using the conditional pdf  $f(y|C_i)$ . Finally, the reconstruction values are taken as the centroid of each quantization cell. The distortion and the average rate of quantizer  $Q_2$  are evaluated using (4) and (11), respectively. To implement the practical FR-SSQ based on the above asymptotic analysis, the same discretization procedure as for the proposed algorithm is utilized with  $K_1 = 3000$  and  $K_2 = 5000$ .

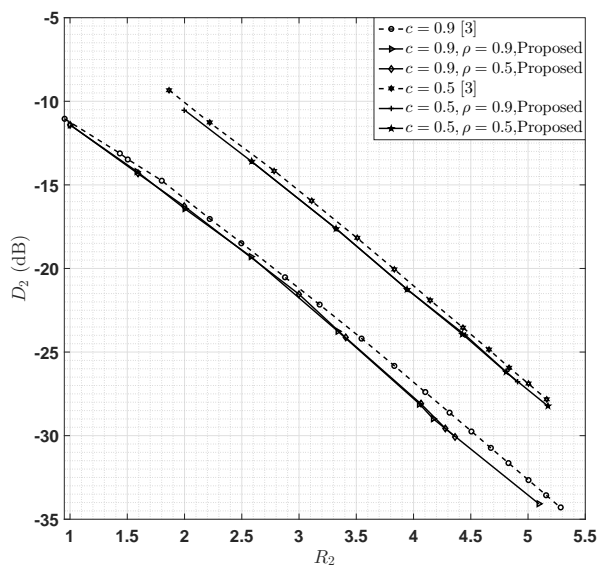
We ran the proposed algorithm for optimal FR-SSQ design for two values of  $M_1$ , namely 4 and 16, for  $\rho = 0.5, 0.9$ , and a set of values of  $\lambda_2$  satisfying  $\lambda_2 \in [0.00001, 0.05]$ .

Figures 8a and 8b plot the distortion  $D_2$  versus the average rate  $R_2$ , for the proposed FR-SSQ in comparison with the scheme of [4], for  $M_1 = 4$  and  $M_1 = 16$ , respectively. The plots for both correlation coefficients  $c = 0.5$  and  $c = 0.9$  and  $\rho = 0.5, 0.9$  are included. It can be observed from both figures that the performance when  $\rho = 0.5$  and  $\rho = 0.9$  is almost identical. It can also be seen that our design always outperforms the scheme of [4]. To make the comparison easier, we show in Tables II and III the performance improvement (in dB) over the scheme of [4] at decoder 2 for various values of  $R_2$ , when  $M_1 = 4$  and 16, respectively. Note that when  $R_2 \approx 1.0$ , the quantizers of  $Y$  for all the schemes have  $M_{2,i} \leq 2$  cells. This explains why the improvement is small at this rate. Then the gap gradually increases with the ascending rates, in most cases. We note that the difference in performance is more pronounced for the higher correlation coefficient and the smaller  $M_1$ . In particular, in the case of  $c = 0.5$ , the improvement is around 0.45 dB for  $2 \leq R_2 \leq 3$ , for both values of  $M_1$ . For  $M_1 = 4$ , the improvement increases as  $R_2$  becomes higher than 3, reaching a peak of 0.75 dB at  $R_2 = 0.47$ , while for  $M_1 = 16$  the performance difference peaks at 0.5 dB. In the case when  $c = 0.9$ , the improvement over the scheme of [4] when  $M_1 = 4$  equals 0.8 dB for  $2 \leq R_2 \leq 3$ , then gradually increases for  $R_2 > 3$ , achieving the value of 1.4 dB when  $R_2 = 4.4$ . For  $M_1 = 16$  the performance gain slightly drops, reaching about 0.55 dB for  $2 \leq R_2 \leq 3$  and a maximum of 1.1 dB at  $R_2 = 4.5$ .

For a fair comparison, we also have to account for the value of  $D_1$ , which is shown in Table IV. We point out that, for fixed  $M_1$ , the value of  $D_1$  obtained with the scheme of



(a)  $M_1 = 4$ .



(b)  $M_1 = 16$ .

Fig. 8: Performance comparison of the proposed FR-SSQ against the level-constrained SSQ of [4].

TABLE II: Performance improvement (in dB) at the second decoder over the scheme of [4] for  $M_1 = 4$ .

$c$ \ $R_2$	1.0	1.5	2.0	2.5	3.0	3.5	4.0	4.4	4.7
0.5	0.003	0.3	0.45	0.45	0.45	0.65	0.55	0.65	0.75
0.9	0.2	0.5	0.8	0.8	0.8	0.95	1.2	1.4	—

[4] is constant, while with our design it varies slightly as  $R_2$  increases up to 3.5, after which it stabilizes. We observe that our scheme outperforms the scheme of [4] at the first decoder when  $M_1 = 4$ , but it is worse when  $M_1 = 16$ . However, the loss in the latter case (which is of only 0.1 dB for  $R_2 \geq 3.5$ ) is offset by the gain in performance at decoder 2. Therefore, we

TABLE III: Performance improvement (in dB) at the second decoder over the scheme of [4] for  $M_1 = 16$ .

$c$ \ $R_2$	1.0	1.5	2.0	2.5	3.0	3.5	4.0	4.5	5.0	5.1
0.5	—	—	0.45	0.45	0.5	0.5	0.5	0.4	0.4	0.4
0.9	0.15	0.4	0.55	0.45	0.55	0.8	1.0	1.1	0.9	0.9

TABLE IV: Comparison of  $D_1$  between the proposed FR-SSQ and the scheme of [4] for  $M_1 = 4$  and 16. The distortion is listed in dB.

$D_1$ \ $M_1$	4		16	
[3]	-9.05		-20.08	
Proposed	$R_2 < 3.5$	$R_2 \geq 3.5$	$R_2 < 3.3$	$R_2 \geq 3.3$
$c = 0.5$	[-9.30, -9.21]	-9.30	[-19.99, -19.89]	-19.99
Proposed	$R_2 < 2.3$	$R_2 \geq 2.3$	$R_2 < 3.5$	$R_2 \geq 3.5$
$c = 0.9$	[-9.30, -9.18]	-9.30	[-20.07, -19.74]	-19.99

conclude that the overall performance of our scheme is higher than that of [4] for both values of  $M_1$ . On the other hand, the performance difference tends to decrease as  $M_1$  increases. This is expected since the asymptotic analysis performed in [4] becomes accurate when the rate approaches infinity.

Figure 9 illustrates the encoder partitions for the proposed FR-SSQ and for the scheme of [4] when  $M_1 = 3$  and  $R_2 \approx 2.17$ . The figure additionally shows the contour of the joint pdf  $f_{XY}$ . It can be noticed that the quantizer of  $X$  for the proposed FR-SSQ (Figure 9a) has more dense outputs in the region where the marginal pdf  $f_X$  takes on large values, compared with the counterpart of [4] (Figure 9b). This could explain the performance improvement of around 0.19 dB in terms of  $D_1$  for our scheme.

It is instructive to examine the probabilities of the cells of the quantizer for  $X$ . For the proposed FR-SSQ, we have  $P(C_1) = 0.2743$ ,  $P(C_2) = 0.4711$  and  $P(C_3) = 0.2546$ , while for the scheme of [4], we have  $P(C_1) = 0.2278$ ,  $P(C_2) = 0.5444$  and  $P(C_3) = 0.2278$ . Note that in both cases,  $P(C_2)$  is higher than  $P(C_1)$  and than  $P(C_3)$ , but cell  $C_2$  is narrower in our design, making its contribution to distortion  $D_1$  smaller than for the scheme of [4]. It turns out that this decrease in the distortion of cell  $C_2$  offsets the resulting increase in the distortion of cells  $C_1$  and  $C_3$ , thus leading to a smaller value of  $D_1$  for our design.

It can also be observed from Figure 9 that in our scheme  $M_{2,1} = M_{2,3} > M_{2,2}$ , while the opposite holds for the design of [4]. This can be attributed to the different constraints imposed in the two designs. Namely, our work constrains the average rate at encoder 2, which is  $\sum_{i=1}^{M_1} P(C_i) \log_2 M_{2,i}$ , to be fixed, while [4] constrains the total number of cells for all encoder 2 quantizers, to be fixed. Since  $P(C_1)$  and  $P(C_3)$  are lower than  $P(C_2)$ , our design allows for values of  $M_{2,1}$  and  $M_{2,3}$  higher than  $M_{2,2}$ , since an extra cell in either  $M_{2,1}$  or  $M_{2,3}$  contributes much less to the average rate than an extra cell in  $M_{2,2}$ . On the other hand, for the design of [4], an extra cell in any quantizer at encoder 2 has the same effect with respect to meeting the constraint. Therefore, more cells are allocated to  $M_{2,2}$  since its distortion has a higher weight in the average distortion  $D_2$  than  $M_{2,1}$  or  $M_{2,3}$ .







It can be easily seen that  $\lim_{K \rightarrow \infty} \|\mathbf{v}_{B,K}^* - \mathbf{v}_B^*\|^2 = 0$ , where  $\|\cdot\|$  denotes the Euclidian norm. Additionally, a moment of thought reveals that  $\beta$  is a continuous function<sup>3</sup>. Thus, we obtain that

$$\mathcal{F}_{EC}(\mathbf{Q}_{\epsilon,B,K}^*, X_B, Y_B) = \mathcal{F}_{EC}(\mathbf{Q}_{\epsilon}^*, X_B, Y_B) + \delta(B, K), \quad (25)$$

for some  $\delta(B, K)$  such that  $\lim_{K \rightarrow \infty} \delta(B, K) = 0$ .

Consider now the EC-SSQ  $\mathbf{Q}_{\epsilon,B,K}^*$  for the pair of discrete RVs  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$ , constructed from  $\mathbf{Q}_{\epsilon,B,K}^*$  as explained next. For each cell  $C_i = (t_u^{(B)}, t_v^{(B)})$  of the first encoder of  $\mathbf{Q}_{\epsilon,B,K}^*$ , the corresponding cell in  $\tilde{\mathbf{Q}}_{\epsilon,B,K}^*$  is  $\tilde{C}_i = \{x_{u+1}^{(B)}, \dots, x_v^{(B)}\}$ . For each cell  $C_{i,j} = (t_m^{(B)}, t_n^{(B)})$  of the second encoder of  $\mathbf{Q}_{\epsilon,B,K}^*$ , the corresponding cell in  $\tilde{\mathbf{Q}}_{\epsilon,B,K}^*$  is  $\tilde{C}_{i,j} = \{y_{m+1}^{(B)}, \dots, y_n^{(B)}\}$ . It follows that  $\mathbb{P}[\tilde{X}_{B,K} \in \tilde{C}_i] = \mathbb{P}[X_B \in C_i]$  and  $\mathbb{P}[\tilde{Y}_{B,K} \in \tilde{C}_{i,j}] = \mathbb{P}[Y_B \in C_{i,j}]$  for all  $i$  and  $j$ . This observation implies that

$$R_{EC,1}(\tilde{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}) = R_{EC,1}(\mathbf{Q}_{\epsilon,B,K}^*, X_B), \quad (26)$$

$$R_{EC,2}(\tilde{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) = R_{EC,2}(\mathbf{Q}_{\epsilon,B,K}^*, X_B, Y_B).$$

Next we will show that the following hold

$$D_1(\tilde{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}) = D_1(\mathbf{Q}_{\epsilon,B,K}^*, X_B) - D_{\tilde{X}_{B,K}}, \quad (27)$$

$$D_2(\tilde{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) = D_2(\mathbf{Q}_{\epsilon,B,K}^*, X_B, Y_B) - \gamma(\mathbf{Q}_{\epsilon,B,K}^*), \quad (28)$$

where

$$D_{\tilde{X}_{B,K}} \triangleq \sum_{k=1}^K \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)})^2 f_{X_B}(x) dx,$$

$$\gamma(\mathbf{Q}_{\epsilon,B,K}^*) \triangleq \sum_{i=1}^{M_1} \sum_{l=1}^K \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} (y - y_l^{(B)})(y + y_l^{(B)} - 2g_2(i, f_2(i, y))) \int_{C_i} f_{X_B Y_B}(x, y) dx dy.$$

Note that  $|y + y_l^{(B)} - 2g_2(i, f_2(i, y))| \leq 4B$  and  $|y - y_l^{(B)}| \leq \frac{2B}{K}$  when  $y \in [t_{l-1}^{(B)}, t_l^{(B)})$ . Thus, we obtain that  $|\gamma(\mathbf{Q}_{\epsilon,B,K}^*)| \leq \frac{8B^2}{K} \sum_{i=1}^{M_1} \sum_{l=1}^K \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} \int_{C_i} f_{X_B Y_B}(x, y) dx dy = \frac{8B^2}{K}$ , which leads to

$$\lim_{K \rightarrow \infty} \gamma(\mathbf{Q}_{\epsilon,B,K}^*) = 0. \quad (29)$$

In order to prove (27) let  $C_i = (t_u^{(B)}, t_v^{(B)})$ . It follows that

$$\begin{aligned} & \int_{t_u^{(B)}}^{t_v^{(B)}} (x - g_1(i))^2 f_{X_B}(x) dx = \\ & \sum_{k=u+1}^v \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)} + x_k^{(B)} - g_1(i))^2 f_{X_B}(x) dx = \\ & \sum_{k=u+1}^v \left( \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)})^2 f_{X_B}(x) dx + \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x_k^{(B)} - g_1(i))^2 f_{X_B}(x) dx + \right. \\ & \quad \left. 2(x_k^{(B)} - g_1(i)) \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)}) f_{X_B}(x) dx \right) = \\ & \sum_{k=u+1}^v \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)})^2 f_{X_B}(x) dx + \sum_{k=u+1}^v (x_k^{(B)} - g_1(i))^2 P_{\tilde{X}_{B,K}}(x_k^{(B)}), \end{aligned}$$

<sup>3</sup>The proof that  $\beta$  is continuous can be found in [27].

where the last equality is due to the fact that  $\int_{t_{k-1}^{(B)}}^{t_k^{(B)}} (x - x_k^{(B)}) f_{X_B}(x) dx = 0$  (based on the definition of  $x_k^{(B)}$ ) and that  $P_{\tilde{X}_{B,K}}(x_k^{(B)}) = \int_{t_{k-1}^{(B)}}^{t_k^{(B)}} f_{X_B}(x) dx$ . The above observation implies (27). Further, in order to prove (28), let  $C_{i,j} = (t_m^{(B)}, t_n^{(B)})$ . It follows that

$$\begin{aligned} & \int_{t_m^{(B)}}^{t_n^{(B)}} (y - g_2(i, j))^2 \int_{C_i} f_{X_B Y_B}(x, y) dx dy \\ & = \sum_{l=m+1}^n \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} \left( (y - y_l^{(B)})(y + y_l^{(B)} - 2g_2(i, j)) \right. \\ & \quad \left. + (y_l^{(B)} - g_2(i, j))^2 \right) \int_{C_i} f_{X_B Y_B}(x, y) dx dy \\ & = \sum_{l=m+1}^n \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} (y - y_l^{(B)})(y + y_l^{(B)} - 2g_2(i, j)) \int_{C_i} f_{X_B Y_B}(x, y) dx dy \\ & \quad + \sum_{l=m+1}^n (y_l^{(B)} - g_2(i, j))^2 \mathbb{P}[\tilde{X}_{B,K} \in \tilde{C}_i, \tilde{Y}_{B,K} = y_l^{(B)}], \end{aligned}$$

where the last equality uses the fact that  $\mathbb{P}[\tilde{X}_{B,K} \in \tilde{C}_i, \tilde{Y}_{B,K} = y_l^{(B)}] = \int_{t_{l-1}^{(B)}}^{t_l^{(B)}} \int_{C_i} f_{X_B Y_B}(x, y) dx dy$ . The above observation implies (28). Further, relations (26)-(28) lead to

$$\mathcal{F}_{EC}(\mathbf{Q}_{\epsilon,B,K}^*, X_B, Y_B) = \mathcal{F}_{EC}(\tilde{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) + \rho D_{\tilde{X}_{B,K}} + (1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*). \quad (30)$$

Further, recall that  $\hat{\mathbf{Q}}_{B,K}$  is the optimal EC-SSQ (with convex cells) for the pair of RVs  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$ . Let  $\mathbf{Q}_{B,K}$  be the corresponding EC-SSQ for  $(X_B, Y_B)$  with thresholds in  $\mathcal{U}_{B,K}$ , according to the correspondence described in the paragraph after equation (25). Then we have, similarly to (30),

$$\mathcal{F}_{EC}(\mathbf{Q}_{B,K}, X_B, Y_B) = \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) + \rho D_{\tilde{X}_{B,K}} + (1 - \rho)\gamma(\mathbf{Q}_{B,K}). \quad (31)$$

Now consider extending the EC-SSQ  $\mathbf{Q}_{B,K}$  to an EC-SSQ  $\bar{\mathbf{Q}}_{B,K}$  for  $X, Y$ , as follows. The encoding partition for  $X$  in  $\bar{\mathbf{Q}}_{B,K}$  has two more cells, namely  $(-\infty, -B)$  and  $(B, \infty)$ , both having the mean of  $X$  as reconstruction. Likewise, when  $X \in [-B, B]$ , the encoder for  $Y$  has two more cells, namely  $(-\infty, -B)$  and  $(B, \infty)$ , both having the mean of  $Y$  as reconstruction. When  $X \notin [-B, B]$ , the encoder for  $Y$  sends only one symbol and the reconstruction is the mean of  $Y$ . It can be readily seen that

$$\mathcal{F}_{EC}(\bar{\mathbf{Q}}_{B,K}, X, Y) = P_B \mathcal{F}_{EC}(\mathbf{Q}_{B,K}, X_B, Y_B) + \epsilon_2(B), \quad (32)$$

for some function  $\epsilon_2(B)$  such that  $\lim_{B \rightarrow \infty} \epsilon_2(B) = 0$ .

The aforementioned discussion implies the following se-

quence of relations

$$\begin{aligned}
 & \mathcal{F}_{EC}^* \\
 & \stackrel{(a)}{\leq} \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, X, Y) \\
 & \stackrel{(b)}{=} P_B \mathcal{F}_{EC}(\mathbf{Q}_{B,K}, X_B, Y_B) + \epsilon_2(B) \\
 & \stackrel{(c)}{=} P_B \left( \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) + \rho D_{\tilde{X}_{B,K}} \right) + \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) + \epsilon_2(B) \\
 & \stackrel{(d)}{\leq} P_B \left( \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{\epsilon,B,K}^*, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) + \rho D_{\tilde{X}_{B,K}} \right) + \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) + \epsilon_2(B) \\
 & \stackrel{(e)}{=} P_B \mathcal{F}_{EC}(\mathbf{Q}_{\epsilon,B,K}^*, X_B, Y_B) + P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) - \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) + \epsilon_2(B) \\
 & \stackrel{(f)}{=} P_B \left( \mathcal{F}_{EC}(\mathbf{Q}_{\epsilon}^*, X_B, Y_B) + \delta(B, K) \right) + P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) - \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) + \epsilon_2(B) \\
 & \stackrel{(g)}{=} \mathcal{F}_{EC}(\mathbf{Q}_{\epsilon}^*, X, Y) - \epsilon_1(B) + P_B\delta(B, K) + \\
 & \quad P_B(1 - \rho)(\gamma(\mathbf{Q}_{B,K}) - \gamma(\mathbf{Q}_{\epsilon,B,K}^*)) + \epsilon_2(B) \\
 & \stackrel{(h)}{\leq} \mathcal{F}_{EC}^* + \epsilon - \epsilon_1(B) + P_B\delta(B, K) + \\
 & \quad P_B(1 - \rho)(\gamma(\mathbf{Q}_{B,K}) - \gamma(\mathbf{Q}_{\epsilon,B,K}^*)) + \epsilon_2(B).
 \end{aligned}$$

Notice that (a) follows from the definition of  $\mathcal{F}_{EC}^*$ , (b) is based on (32), (c) follows from (31), (d) holds in virtue of the optimality of  $\hat{\mathbf{Q}}_{B,K}$  for  $(\tilde{X}_{B,K}, \tilde{Y}_{B,K})$ , (e) follows from (30), (f) from (25), (g) is based on (24) and (h) is based on (23). Next we use the sequence of relations (a)–(h) and apply the fact that  $A_1 \leq A_2 \leq A_3$  implies that  $A_2 - A_1 \leq A_3 - A_1$  and  $A_3 - A_2 \leq A_3 - A_1$ , for  $A_1 = \mathcal{F}_{EC}^*$ ,  $A_2$  being the right hand side of (c) and  $A_3$  being the right hand side of (h). Thus, we obtain

$$\begin{aligned}
 & P_B \left( \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) + \rho D_{\tilde{X}_{B,K}} \right) + \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) + \epsilon_2(B) - \mathcal{F}_{EC}^* \leq \\
 & \epsilon - \epsilon_1(B) + P_B\delta(B, K) + P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) - \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) + \epsilon_2(B).
 \end{aligned} \tag{33}$$

$$\begin{aligned}
 & \mathcal{F}_{EC}^* + \epsilon - \epsilon_1(B) + P_B\delta(B, K) + P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) - \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) + \epsilon_2(B) - \\
 & \quad P_B \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) - \\
 & \quad P_B \left( \rho D_{\tilde{X}_{B,K}} + (1 - \rho)\gamma(\mathbf{Q}_{B,K}) \right) - \epsilon_2(B) \leq \\
 & \epsilon - \epsilon_1(B) + P_B\delta(B, K) + P_B(1 - \rho)\gamma(\mathbf{Q}_{B,K}) - \\
 & \quad P_B(1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) + \epsilon_2(B).
 \end{aligned} \tag{34}$$

Relation (33) implies that

$$\begin{aligned}
 & \mathcal{F}_{EC}^* - P_B \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) \geq \\
 & \quad P_B \left( \rho D_{\tilde{X}_{B,K}} - \delta(B, K) + (1 - \rho)\gamma(\mathbf{Q}_{\epsilon,B,K}^*) \right) + \epsilon_1(B) - \epsilon.
 \end{aligned}$$

Relation (34) leads to

$$\begin{aligned}
 & \mathcal{F}_{EC}^* - P_B \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) \leq \\
 & \quad P_B \left( \rho D_{\tilde{X}_{B,K}} + (1 - \rho)\gamma(\mathbf{Q}_{B,K}) \right) + \epsilon_2(B).
 \end{aligned}$$

The above two inequalities, together with (29) and  $\lim_{K \rightarrow \infty} D_{\tilde{X}_{B,K}} = \lim_{K \rightarrow \infty} \gamma(\mathbf{Q}_{B,K}) = \lim_{K \rightarrow \infty} \delta(B, K) = 0$ ,  $\lim_{B \rightarrow \infty} \epsilon_1(B) = \lim_{B \rightarrow \infty} \epsilon_2(B) = 0$  and  $\lim_{B \rightarrow \infty} P_B = 1$ , lead to

$$0 \leq \lim_{B \rightarrow \infty} \lim_{K \rightarrow \infty} \mathcal{F}_{EC}(\hat{\mathbf{Q}}_{B,K}, \tilde{X}_{B,K}, \tilde{Y}_{B,K}) - \mathcal{F}_{EC}^* \leq \epsilon,$$

for every  $\epsilon > 0$ , which implies that relation (21) holds. Thus, the proof is completed. ■

#### ACKNOWLEDGMENT

The authors would like to thank the Editor and the anonymous reviewers for their valuable comments and suggestions, which helped improve the quality of the work.

#### REFERENCES

- [1] H. Viswanathan and T. Berger, "Sequential coding of correlated sources," *IEEE Trans. Inform. Theory*, vol. 46, no. 1, pp. 236-246, Jan. 2000.
- [2] J. Wang, X. Wu, J. Sun and S. Yu, "On two-stage sequential coding of correlated sources," *IEEE Trans. Inform. Theory*, vol. 60, no. 12, pp. 7490-7505, Dec. 2014.
- [3] B. Rimoldi, "Successive refinement of information: characterization of achievable rates," *IEEE Trans. Inform. Theory*, vol. 40, no. 1, pp. 253-259, Jan. 1994.
- [4] R. Balasubramanian, C. A. Bouman, and J. P. Allebach, "Sequential scalar quantization of vectors: an analysis," *IEEE Trans. Image Proc.*, vol. 4, no. 9, pp. 1282-1295, Sept. 1995.
- [5] J. Z. Chang and J. P. Allebach, "Optimal sequential scalar quantization of vectors", in *Proc. 27th Asilomar Conf. Signal, Systems and Computers*, pp. 966-971, Pacific Grove, CA, USA, Nov. 1993.
- [6] R. Balasubramanian, C. A. Bouman, and J. P. Allebach, "Sequential scalar quantization of color images", *J. Electron. Imaging.*, vol. 3, no. 1, pp. 45-59, Jan. 1994.
- [7] S. P. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inform. Theory*, vol. IT-28, no. 2, pp. 129-137, Mar. 1982.
- [8] A. V. Trushkin, "Sufficient conditions for uniqueness of a locally optimal quantizer for a class of convex error weighting functions", *IEEE Trans. Inform. Theory*, vol. 28, no. 2, pp. 187-198, Mar. 1982.
- [9] J. C. Kieffer, "Uniqueness of locally optimal quantizer for log-concave density and convex error weighting function", *IEEE Trans. Inform. Theory*, vol. IT-29, no. 1, pp. 42-47, Jan. 1983.
- [10] S. Dumitrescu and X. Wu, "On properties of locally optimal multiple description scalar quantizers with convex cells," *IEEE Trans. Inform. Theory*, vol. 55, no. 12, pp. 5591-5606, Dec. 2009.
- [11] S. Dumitrescu, "On the design of optimal noisy channel scalar quantizer with random index assignment," *IEEE Trans. Inform. Theory*, vol. 62, no. 2, pp. 724-735, Feb. 2016.
- [12] A. Gyorgy and T. Linder, "On the structure of optimal entropy-constrained scalar quantizers", *IEEE Trans. Inform. Theory*, vol. 48, no. 2, pp. 416-427, Feb. 2002.
- [13] D. Muresan and M. Effros, "Quantization as histogram segmentation: optimal scalar quantizer design in network systems," *IEEE Trans. Inform. Theory*, vol. 54, no. 1, pp. 344-366, Jan. 2008.
- [14] D. K. Sharma, "Design of absolutely optimal quantizers for a wide class of distortion measures", *IEEE Trans. Inform. Theory*, vol. 24, no. 6, pp. 693-702, Nov. 1978.
- [15] X. Wu, "Optimal quantization by matrix searching", *J. Algorithms*, vol. 12, no. 4, pp. 663-673, Dec. 1991.
- [16] X. Wu and K. Zhang, "Quantizer monotonicities and globally optimal scalar quantizer design," *IEEE Trans. Inform. Theory*, vol. 39, no. 3, pp. 1049-1053, May. 1993.
- [17] S. Dumitrescu and X. Wu, "Optimal multiresolution quantization for scalable multimedia coding," in *Proc. IEEE Information Theory Workshop (ITW 2002)*, pp. 139-142, Bangalore, India, Oct. 2002.
- [18] S. Dumitrescu and X. Wu, "Algorithms for optimal multi-resolution quantization," *J. Algorithms*, vol. 50, no. 1, pp. 1-22, Jan. 2004.
- [19] S. Dumitrescu and X. Wu, "Fast algorithms for optimal two-description scalar quantizer design", *Algorithmica*, vol. 41, no. 4, pp. 269-287, Feb. 2005.
- [20] S. Dumitrescu, X. Wu, "Lagrangian optimization of two-description scalar quantizers," *IEEE Trans. Inform. Theory*, vol. 53, no. 11, pp. 3990-4012, Nov. 2007.

- [21] A. Gyorgy, T. Linder, G. Lugosi, "Tracking the best quantizer" *IEEE Trans. Inform. Theory*, vol. 54, no. 4, pp. 1604-1625, Apr. 2008.
- [22] H. Wu and S. Dumitrescu, "Design of optimal entropy-constrained scalar quantizer for sequential coding of correlated sources", in *Proc. IEEE Inform. Theory Workshop (ITW 2017)*, pp. 524-528, Kaohsiung, Taiwan, Nov. 2017.
- [23] P. A. Chou, T. Lookabaugh, and R. M. Gray, "Entropy-constrained vector quantization", *IEEE Trans. Acoust., Speech, Signal Proc.*, vol. 37, no. 1, pp. 31-42, Jan. 1989.
- [24] M. Fleming, Q. Zhao, and M. Effros, "Network vector quantization," *IEEE Trans. Inform. Theory*, vol. 50, no. 8, pp. 1584-1604, Aug. 2004.
- [25] A. Aggarval, M. Klave, S. Moran, P. Shor, and R. Wilber, "Geometric applications of a matrix-searching algorithm", *Algorithmica*, vol. 2, no. 1-4, pp.195-208, Nov. 1987.
- [26] H. Gish and N. J. Pierce, "Asymptotically efficient quantizing," *IEEE Trans. Inform. Theory*, vol. IT-14, no. 5, pp. 676-683, Sep. 1968.
- [27] H. Wu, "Sequential scalar quantization of two dimensional vectors in polar and cartesian coordinates," *Ph.D. Dissertation*, McMaster University, Hamilton, ON, Canada, Aug. 2018.
- [28] A. Gyorgy, T. Linder, P. A. Chou, and B. J. Betts, "Do optimal entropy-constrained quantizers have finite or infinite number of codewords?", *IEEE Trans. Inform. Theory*, vol. 49, no. 11, pp. 3031-3037, Nov. 2003.

PLACE  
PHOTO  
HERE

**Huihui Wu** (S'14) received the B.Sc. degree in communication engineering from Southwest University for Nationalities, Chengdu, China, in 2011, and the M.S. degree in communication engineering from Xiamen University, Xiamen, China, in 2014. He received the Ph.D. degree in electrical and computer engineering from McMaster University, Hamilton, Canada. His research interests include channel coding, joint source and channel coding, multiple description coding, and signal quantization.

PLACE  
PHOTO  
HERE

**Sorina Dumitrescu** (M'05-SM'13) received the B.Sc. and Ph.D. degrees in mathematics from the University of Bucharest, Romania, in 1990 and 1997, respectively. From 2000 to 2002 she was a Postdoctoral Fellow in the Department of Computer Science at the University of Western Ontario, London, Canada. Since 2002 she has been with the Department of Electrical and Computer Engineering at McMaster University, Hamilton, Canada, where she held a Postdoctoral and a Research Associate position, and where she is currently an Associate

Professor. Her current research interests include multimedia coding and communications, network-aware data compression, joint source-channel coding, signal quantization. Her earlier research interests were in formal languages and automata theory. She was a recipient of the NSERC University Faculty Award during 2007-2012.