# HIGH-RATE DESIGN OF TRANSFORM CODERS WITH GAUSSIAN MIXTURE COMPANDERS

Ethan R. Duni and Bhaskar D. Rao

Department of Electrical and Computer Engineering University of California, San Diego La Jolla, CA 92093-0407 Email: {eduni,brao}@ucsd.edu

#### ABSTRACT

This paper examines the problem of designing fixed-rate transform coders for sources with arbitrary distributions, under input-weighted squared error distortion measures. As a component of this system, a flexible scalar compander using Gaussian Mixtures is proposed. An algorithm is developed to set the parameters of the system using a data-driven technique that automatically balances the source statistics, distortion measure, and structure of the transform coder to minimize the high-rate distortion. The implementation of Gaussian Mixture companders is explored, resulting in a flexible, lowcomplexity scalar quantizer. The operation of this system for the problem of wideband speech spectrum quantization with Log Spectral Distortion is illustrated, and shown to provide good performance with very low, rate-independent complexity.

## 1. INTRODUCTION

Transform coding is a popular method for quantizing vectors of data using only scalar quantizers. Figure 1 illustrates the structure of a transform coder using companding scalar quantizers. As a result of this structure, transform coders have a number of desirable properties, such as small storage requirements and very low coding complexity. The price for these features is inferior performance as compared to more flexible quantizers.



Fig. 1. Transform coder with compandering scalar quantizers

Many approaches to transform coding depend on fixing the transform in some way, often by assuming the source is Gaussian (see [1],[2]). In other cases, only certain convenient transforms are considered: an example would be the DCT in image coding [3]. This paper considers a more general problem in which the distribution is unknown. Additionally, this paper considers more general distortion measures, such as Log Spectral Distortion. To accommodate these concerns, an algorithm is developed to set the parameters of the system using a data-driven design technique that automatically balances the source statistics, distortion measure, and structure of the transform coder to achieve minimal high rate distortion. To allow for unknown data statistics, a flexible compander system using Gaussian Mixtures is proposed.

The problem of designing a transform coder for minimum distortion from a database has also been considered by Archer and Leen in [4]. This paper focuses on fixed-rate systems and considers a more general distortion measure, whereas [4] covers the variable-rate case with MSE distortion. Also, this training method is rate-independent, which allows operation at arbitrary rates with no additional storage or training requirements. Another similar work is [3], which considers the variable-rate transform coding of images and uses Gaussian mixtures to model the marginal source densities. In that work, the transform is considered as fixed and only the problem of learning the component scalar quantizers is addressed.

To illustrate the performance of the proposed system, this paper will focus on the example of wideband speech LSF quantization, with the Log Spectral Distortion measure. An LSF quantization system is said to achieve transparent quality when its average LSD is no more than 1dB, produces outliers between 2-4dB less than 1% of the time and produces a negligible amount of larger outliers. In [5], Gosset lattices are utilized for wideband speech LSF quantization, attaining transparent quality at around 45 bits per frame. In [6], MSVQ is employed and results in transparent quality at around 56 bits per frame. In [7], a Gaussian Mixture Model based VQ system is shown to achieve transparent quality at 48 bits per frame. All of these works consider only memoryless systems. It is noteworthy that essentially all authors have found that rates of at least 45 bits per vector are required for high quality memoryless quantization of wideband LSF vectors. For such large codebooks, the transform coder offers extremely low complexity, as one need only store the system parameters. Furthermore, if the scalar quantizers are implemented with companders, the encoding complexity becomes rateindependent. Although these features prevent the transform coder from performing as well as the best schemes, it is still able to improve upon the performance of MSVQ. Thus, the transform coder is attractive in scenarios where very low complexity is required, and so a framework is developed for achieving the best possible high rate performance.

This research was supported by Micro Grants 03-073, 04-074 and 05-033, sponsored by Qualcomm Inc.

## 2. BACKGROUND

Consider a source  $X \in \mathbb{R}^d$ , with density f(x). This paper will restrict attention to *input-weighted* distortion measures of the form

$$d(x_1, x_2) = (x_1 - x_2)^{\mathsf{T}} S(x_1)(x_1 - x_2)$$
(1)

were S(x) is a symmetric, positive-definite matrix, called the *sensitivity matrix*. Analysis of fixed-rate quantization under inputweighted distortion measures can be found in [8]. A wide variety of distortion measures, including Log Spectral Distortion, can be accurately approximated at high rates in this fashion. Calculation of the sensitivity matrix for LSD on LSF vectors is detailed in [8].

This paper will restrict attention to the high-rate case, where the distortion of an *r*-bit quantizer can be approximated as an integral, as in [1]. In the transform coder, all cells are hyperrectangles. Due to its structure, the transform coder suffers from space-filling loss, oblongitis, and a limited ability to exploit dependence between the elements of X (see [9] for a detailed explanation). The high-rate approximation for a transform coder under MSE is given in [1]. Let  $\lambda_{\theta_i}(y_i)$  be the point densities of the scalar quantizers, and  $K_i$  the numbers of levels assigned to each of them. For an *r*-bit quantizer,  $K_i$  is parameterized in terms of the new variables  $\beta_i$  as follows:

$$K_i = 2^{r/d} \beta_i \left( \prod_{j=1}^d \beta_j \right)^{-1/d}$$

Here,  $\beta_i > 0$ , insuring that  $K_i > 0$ . The form of this parametrization insures that  $\prod_{i=1}^{d} K_i = 2^r$ . For training purposes, the constraint that the  $K_i$ 's must be integers is ignored. Instead, a pruning algorithm is applied to meet the constraint when implementing the coder. The high-rate distortion of an *r*-bit transform coder is [11]:

$$D_{\theta} \cong 2^{-2r/d} \frac{1}{12} \left( \prod_{j=1}^{d} \beta_{j} \right)^{2/d} \sum_{i=1}^{d} \beta_{i}^{-2} \operatorname{E} \left( ||t_{i}||_{S(X)}^{2} \lambda_{\theta_{i}}^{-2}(t_{i}^{\mathsf{T}}X) \right)$$
(2)

where  $t_i$  is the *i*-th column of *T*.

### 2.1. Point Densities

Next, a specific parametric form for the compander point densities is required. This paper proposes a flexible class of point densities which are Gaussian Mixtures:

$$\lambda_{\theta_i}(y_i) = \sum_{m=1}^M \alpha_{im} N(y_i | \mu_{im}, \sigma_{im}^2)$$

This class of point densities can, as M grows large, approximate a wide variety of densities. Even for low values of M, mixtures are able to model features such as multimodality and skew. Thus, the system parameters are  $\theta = \{T, \beta, \alpha, \mu, \sigma^2\}$  with the constraints  $T^T T = I, \beta_i > 0, \sum_{m=1}^{M} \alpha_{im} = 1$  and  $\sigma_{im}^2 > 0$ . The remainder of this paper will use the notation  $N_{im}(y_i) = N(y_i | \mu_{im}, \sigma_{im}^2)$ .

#### 3. DATA-DRIVEN TRANSFORM CODER DESIGN

The classic approach to designing transform coders is to assume the source is Gaussian. Then, the parameter settings that minimize the high rate distortion can be derived (see [1]). Specifically, T is the KLT,  $\beta_i$  is the square root of *i*-th eigenvalue of the source covariance, and  $\lambda_i(y_i)$  is a Gaussian with variance three times the *i*-th

eigenvalue. As discussed in [10], however, this design can be very poor if the source is not actually Gaussian, and so this paper proposes a numerical technique for designing the system. In practice, one has no knowledge of the distribution except for a set of samples  $x_1, \ldots, x_N$ , drawn i.i.d. from f(x). Then, the expectations in Eq. (2) are replaced by averages, and the design problem is:

$$\min_{\theta \in \Theta} \left( \prod_{j=1}^{d} \beta_j \right)^{2/d} \sum_{i=1}^{d} \beta_i^{-2} \sum_{n=1}^{N} ||t_i||_{S(x_n)}^2 \left( \sum_{m=1}^{M} \alpha_{im} N_{im}(t_i^{\mathsf{T}} x_n) \right)^{-2}$$
(3)

This design procedure is intended to be carried out off-line, and so its complexity does not come to bear on the operation of the resulting transform coder. Also, while successive LSF frames are actually correlated, this dependence is ignored here, as all of the systems under consideration are memoryless.

It is difficult to optimize Eq. (3) over all parts of  $\theta$  simultaneously. As such, an iterative algorithm that alternatingly optimizes over subsets of parameters while holding the others fixed is used. Specifically, each iteration first optimizes over the transform, then over the point density parameters, and then over the level allocations. Full details of this scheme can be found in [11], and are summarized here with an emphasis on the point density optimization step, which is the most novel part. To optimize the transform, this work utilizes the steepest descent approach proposed by Manton in [12]. Notably, it requires only the evaluation of the derivative of Eq. (3) with respect to T, and uses an SVD to enforce the orthogonality constraint. A linesearch approach is used to ensure convergence. The point density parameters are optimized with an extension of the EM algorithm, which is detailed in Section 3.1. Finally, the level allocation problem can be solved by taking the logarithm of Eq. (3), taking its derivative with respect to  $\beta_i$  and setting it to zero, giving:

$$\frac{\beta_i^{-2}c_i}{\sum_{j=1}^d \beta_j^{-2} c_j} = \frac{1}{d}$$

In other words, one should choose the  $\beta_i$ 's such that  $\beta_1^{-2}c_1 = \beta_2^{-2}c_2 = \ldots = \beta_d^{-2}c_d$ . This can be easily accomplished by setting  $\beta_1 = 1$  and then applying the equation  $\beta_i = \sqrt{\frac{c_i}{c_1}}$ ,  $\forall i \ge 2$ .

To initialize the algorithm, one needs an initial guess of the transform. Two obvious choices are the KLT and identity matrices, which will always be included in this work. Given an initial transform, the K-means algorithm is applied to each dimension of the transformed data to initialize the point densities. The level allocation are then initialized as above. Further, it is often useful to perform the training in a hierarchical manner. That is, first perform the training for M = 1, which has much less sensitivity to local minima, then utilize the resulting transform to initialize for M = 2, and so on with the higher orders.

#### 3.1. Point Density Optimization

Notice that the overall objective function, Eq. (3), is a sum over functions of the different scalar quantizers, and so they may be optimized independently. For the *i*-th quantizer, the optimization problem is:

$$\min_{\theta_i} \sum_{n=1}^{N} ||t_i||_{S(x_n)}^2 \left( \sum_{m=1}^{M} \alpha_{im} N_{im}(t_i^{\mathsf{T}} x_n) \right)^{-2} \tag{4}$$

Such a problem can be approached with an extension of the EM algorithm. Where conventional EM applies Jensen's inequality to a logarithm to construct a bound on the objective function (see [13]),

the same can be done with the power function  $(\cdot)^{-2}$ , resulting in the following problem:

$$\min_{\theta_i} \sum_{n=1}^{N} ||t_i||_{S(x_n)}^2 \sum_{m=1}^{M} r_{mn}^3 \left( \alpha_{im} N_{im}(t_i^{\mathsf{T}} x_n) \right)^{-2}$$
(5)

Eq. (5) is an upper bound on Eq. (4). The setting of  $r_{mn}$  is the same as in conventional EM:

$$r_{mn} = \frac{\alpha_{im} N_{im}(t_i^{\mathsf{T}} x_n)}{\sum_{m=1}^{M} \alpha_{im} N_{im}(t_i^{\mathsf{T}} x_n)} \tag{6}$$

Thus, one may construct an iterative optimization procedure for  $\theta_i$  by starting with some initial guess  $\theta_i^0$  and then alternating between an E-step (Eq. (6)) and an M-step that optimizes the resulting bound. Optimization of Eq. (5) over the mixture weights is accomplished by the standard Lagrange multiplier approach, giving:

$$\alpha_{im} \propto \left(\sigma_{im}^2 \sum_{n=1}^N ||t_i||_{S(x_n)}^2 r_{mn}^3 e^{\sigma_{im}^{-2}(x_n - \mu_m)^2}\right)^{1/3}$$
(7)

Optimization of the means and variances of each component is not possible in closed form. In practice, Newton's method works fine on both problems, provided a reasonable initializer is used. Complete details of the Newton steps for the means and variances can be found in [11]. A simple linesearch scheme is used to ensure that the resulting variance estimates are positive.

## 4. IMPLEMENTATION OF GMM TRANSFORM CODER

This section discusses the implementation of a transform coder with GMM point densities. To operate such a system, one first must compute the level allocation for the desired rate. A pruning algorithm is used to make the resulting total number of codepoints as close to  $2^r$  as possible without exceeding it, as described in [11].

The other implementation issue is the evaluation of the compander functions  $g(y_i)$  and  $h(z_i) = g^{-1}(z_i)$ . For the case of a Gaussian Mixture density, the compressor function is easy to evaluate:  $g(y_i) = \sum_{m=1}^{M} \alpha_{im} \Phi_{im}(y_i)$ , where  $\Phi_{im}(y_i)$  denotes the cdf of the *im*-th Gaussian. This expression is easy to evaluate numerically, requiring only routines for computing the error function. However, the expander function is difficult to evaluate, since one cannot interchange the inverse with the summation.

To get around the difficulty with the expander function, Newton's root-finding method can be used to compute  $h(z_i)$ . This procedure begins with an initial guess  $y_i^0$  and then performs a number of iterations to improve this estimate. As is well known, Newton's method provides quadratic convergence if one supplies a suitable initializer. Thus, a simple method for selecting a good initializer for any possible value of  $z_i$  is needed. This can be accomplished by partitioning  $g(y_i)$  into concave and convex regions, and assigning a different initializer depending on which region  $z_i$  falls into. The method of partitioning is shown in Figure 2. The Newton iteration for this problem is:

$$y_i^{k+1} = y_i^k - \frac{\sum_{m=1}^M \alpha_{im} \Phi_{im}(y_i^k) - z_i}{\sum_{m=1}^M \alpha_{im} N_{im}(y_i^k)}$$
(8)

For speech signals, it has been observed that ten iterations of this algorithm are sufficient to result in an average decoding error on the order of  $10^{-32}$ , with a maximum decoding error on the order



Fig. 2. Illustration of the expander initialization scheme, showing (a) GM CDF and (b) its second derivative. Given some  $z \in [0, 1]$  to decode, the initializer is formed by finding the partition that contains z (the dashed lines) and using the inflection point y within that region (the dotted lines).

of  $10^{-31}$ . This is sufficient accuracy for rates up to approximately 30 bits per dimension, which is extremely high. For more moderate rates, five or so iterations are sufficient.

#### 5. PRACTICAL RESULTS

The problem of wideband speech spectrum coding, under the Log Spectral Distortion measure, is considered. It should be noted that, for the high-rate analysis, LSD is measured in dB<sup>2</sup> in order to correspond to input-weighted squared error, and so high-rate approximations use this metric (Figure 3, specifically). However, when calculating operating point and outlier statistics (as in Table 1), the more conventional approach of measuring in dB is used. A training set of 300,000 wideband speech LSF vectors of order 16 was gathered, and the LSD-sensitivity matrix was evaluated for each vector. For testing purposes, an independent database of 65,000 LSF vectors was employed. First, a transform coder was designed using single Gaussian scalar quantizers. To initialize the parameters, the data was assumed to be Gaussian and so the sample mean and covariance were used to set the parameters to be optimal under for Gaussian, MSE assumption. The data-driven transform coder design algorithm described in Section 3 was then used to optimize over the actual statistics and distortion measure. After trying a variety of other initializers (T = I), and many random transforms), it was determined that this result was indeed the best possible.

The performance before and after optimization, both actual and predicted, is seen in Figure 3. Notice that, at high rates, optimization resulted in a large savings of around 8 bits per dimension (128 bits per frame). However, these results are only valid at very high rates, above 20 bits per dimension. At very low rates, the non-optimized



**Fig. 3.** High rate performance. Experimental results in bold lines, high-rate predictions in dashed lines.



**Fig. 4.** Histograms and point densities for transform coefficients. Histograms in solid lines, corresponding optimal point densities in dashed lines and actual point densities in dotted lines.

system performs slightly better, as the optimized system has changed  $\mu_{im}$  to compensate for skew. While the average performance of the systems around the operating point are similar, the optimized system produces far fewer outliers.

Next, the number of Gaussians used in each scalar quantizer was increased. No significant changes in the high rate distortion were observed for M > 2, and so a mixture of 2 Gaussians are used for each compander. This resulted in a reduction in high-rate distortion of 10% over the M = 1 case. The performance of this system, including outlier statistics, is listed in Table 1, where it is compared to the initial, Gaussian-assumption design. The average distortion of each system is roughly the same at rates of interest. However, the outlier statistics have been greatly improved by utilizing GM companders and high-rate optimization. The point densities of two of the

bits/frame	Avg. LSD	Outliers (in %)	
	(in dB)	2-4 dB	> 4  dB
53	1.1094	1.3417	0.0046
	1.0589	1.3509	0.0031
54	1.0175	0.8133	0.0031
	1.0197	1.1045	0.0046
55	0.9804	0.6501	0.0046
	0.9808	0.8426	0.0046
56	0.9419	0.5376	0.0046
	0.9645	0.7969	0.0031

**Table 1.** Performance Around 1dB LSD Operating Point. For each rate, the upper row corresponds to an optimized coder with M = 2, and the lower to a system designed under the Gaussian assumption.

optimized companders are seen in Figure 4, along with histograms of the transform coefficients. The GM compander is able to account for multimodality and skew, and so closely approximates the optimal point densities. Overall, the system outperforms MSVQ by 1-2 bits per frame, with very low, rate independent complexity. While more complex systems offer better performance, transform coding is an attractive option when very low complexity is desired.

# 6. REFERENCES

- [1] S. Na and D. L. Neuhoff, "Bennett's Integral for Vector Quantizers", IEEE Trans. on Info. Theory, Vol. 41, (no.4), July 1995.
- [2] J.J.Y. Huang and P.M. Schultheiss, "Block Quantization of Correlated Gaussian Random Variables" IEEE Trans. on Comm. Systems, Vol. 11, Issue 3, Sept. 1963.
- [3] J. K. Su and R. M. Mersereau, "Coding Using Gaussian Mixture and Generalized Gaussian Models" International Conf. on Image Proc., Sept. 1996.
- [4] C. Archer and T. K. Leen, "A Generalized Lloyd-Type Algorithm for Adaptive Transform Coder Design" IEEE Trans. on Signal Proc., Vol. 52, (no.1), January 2004.
- [5] S. Ragot, J.-P. Adoul, R. Lefebvre and R. Salami, "Low Complexity LSF Quantization for Wideband Speech Coding" IEEE Workshop on Speech Coding, 1999.
- [6] M. Ferhaoui and S. Van Gerven, "LSP Quantization In Wideband Speech Coders" IEEE Workshop on Speech Coding, 1999.
- [7] A. D. Subramaniam, W. R. Gardner and B. D. Rao, "Low Complexity Source Coding Using Gaussian Mixture Models, Lattice Vector Quantization, and Recursive Coding with Application to Speech Spectrum Quantization" IEEE Trans. on Speech and Audio Proc., 2006.
- [8] W. R. Gardner and B. D. Rao, "Theoretical Analysis of the High-Rate Vector Quantization of LPC parameters", IEEE Trans. on Speech and Audio Proc., Vol.3, (no.5), September 1995.
- [9] T. D. Lookabaugh and R. M. Gray, "High-Resolution Quantization Theory and the Vector Quantizer Advantage" IEEE Trans. on Info. Theory, Vol. 35, (no.5), September 1989.
- [10] M. Effros, H. Feng and K. Zeger, "Suboptimality of the Karhunen-Love Transform for Transform Coding" IEEE Trans. on Info. Theory, Vol. 50, (no.8), August 2004.
- [11] E. R. Duni and B. D. Rao, "A High Rate Optimal Transform Coder with GMM Companders" Submitted to IEEE Trans. on Speech and Audio Proc., Sept. 2005.
- [12] J. H. Manton, "Optimization Algorithms Exploiting Unitary Constraints" IEEE Trans. on Signal Proc., Vol. 50, (no.3), March 2002.
- [13] R. M. Neal and G. E. Hinton, "A View of the EM Algorithm that Justifies Incremental, Sparse, and Other Variants" Learning in Graphical Model, M.I. Jordan (editor).