

Robust Quantization of Memoryless Sources Using Dispersive FIR Filters

Kris Popat and Kenneth Zeger

Abstract—A novel approach to quantizing discrete-time memoryless sources is presented. An important feature is that its performance is largely insensitive to errors in modeling the input pdf. The method involves changing the amplitude distribution of the source to be approximately Gaussian by all-pass filtering, then applying a Lloyd–Max quantizer designed for a Gaussian source. After quantization, the samples are passed through another all-pass filter, which is an approximate inverse of the first filter. The mean-square error (MSE) for the overall process is roughly equal to the quantization MSE for the intermediate Gaussian signal, independent of the source statistics. For some sources, this is actually an improvement over direct, correct-model Lloyd–Max quantization. The cost of this technique is some delay due to filtering.

I. INTRODUCTION

MUCH has been written about quantization of memoryless sources, in particular, Laplacian and gamma sources [1], [2]. The subject is important because these sources are often used as models in image and speech coding [3], [4].

An irony associated with the quantization of Laplacian and gamma sources is made evident by the graphs in Fig. 1. Although the rate-distortion functions of these sources are quite promising relative to, say, that of a Gaussian source, simple quantization¹ does not fulfill that promise. In fact, it can be seen that for any given rate, the Lloyd–Max quantizer achieves a lower MSE for Gaussian sources than for Laplacian or gamma sources. Errors in modeling the source further reduce the performance of Lloyd–Max quantization. These observations lead to the quantization scheme described in this paper.

The present suggestion is to use simple quantization, but to filter the source before and after quantizing. If the filters are appropriately designed (see Section IV), then the filtered input signal will have an approximately Gaussian distribution, and the resulting quantization MSE of the overall system will approximate that for direct quantization of a Gaussian source. Moreover, since the initial filtering will tend to make any memoryless source appear Gaussian, the performance of the system is insensitive to errors in modeling the input. This

Paper approved by the Editor for Speech Processing of the IEEE Communications Society. Manuscript received March 8, 1990; revised October 15, 1991 and January 15, 1992. This work was supported in part by Hewlett-Packard Co. and the National Science Foundation. This paper was presented in part at the 1990 International Symposium on Information Theory and its Applications (ISITA '90), Honolulu, HI, November 1990.

K. Popat is with M.I.T. Media Laboratory, Cambridge, MA 02139.

K. Zeger is with the Coordinated Science Laboratory, University of Illinois, Urbana, IL 61801.

IEEE Log Number 9203571.

¹Throughout, the term *simple quantization* refers to fixed-rate minimum MSE memoryless scalar quantization [5], [6].

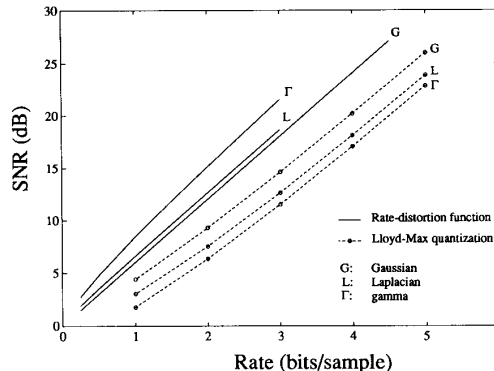


Fig. 1. Performance (mean-square error versus rate) of Lloyd–Max quantization of Laplacian, gamma, and Gaussian sources, relative to the respective rate-distortion functions. Samples of the rate-distortion function for Laplacian and gamma sources were computed by means of the Blahut algorithm [16]. Quantizer performance figures are taken from [3, p. 135].

robustness to the source statistics is a valuable feature, not normally present in quantization systems.

Throughout this paper, it is assumed that the source is stationary and memoryless. For simplicity of notation, it is further assumed that the source has zero-mean and unit-variance.

Many sophisticated alternatives to simple quantization have been suggested for memoryless sources. These alternatives include vector quantization [2] entropy-coded (variable-rate) quantization [7] and trellis coded quantization [8], [9]. While these techniques generally achieve better signal-to-noise ratios than the proposed scheme, they are often more complex, and do not assure the same degree of robustness.

The method of quantization proposed here appears to be novel, despite its stark simplicity. The only similar proposal of which the authors are aware is one by Strube [10]. In that scheme, an all-pass filter is used in a speech ADPCM system to disperse pitch pulses over time, so that quantizer overload-distortion is reduced. Strube's scheme does *not* make use of an inverse filter, however, so that it does not result in a signal that approximates the original. The use of an all-pass prefilter and inverse postfilter to reversibly change the PDF of a signal has been suggested by Zenith [11] in the context of transmission of high-definition television signals. However, their suggestion has nothing to do with quantization.

II. PRESERVATION OF QUANTIZATION MSE

In this section it is shown that the MSE for the proposed system is nearly equal to that incurred by quantizing the signal after prefiltering.

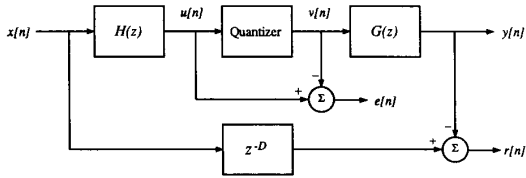


Fig. 2. Quantization system employing a prefilter and a postfilter.

Let $H(z)$ and $G(z)$ denote the z -transforms of the prefilter and postfilter, respectively, as shown in Fig. 2, and let $h[n]$ and $g[n]$ denote the corresponding impulse responses. Let

$$r[n] = x[n - D] - y[n] \quad (1)$$

denote the error of the overall system where D is the delay due to filtering, and let $e[n]$ denote the error incurred by quantizing the prefiltered signal $u[n]$ into $v[n]$. In the following analysis, it is *not* assumed that $e[n]$ is independent of $u[n]$.

By the linearity of the postfilter,

$$\begin{aligned} y[n] &= g[n] * (u[n] - e[n]) \\ &= g[n] * u[n] - g[n] * e[n] \\ &= g[n] * h[n] * x[n] - g[n] * e[n] \end{aligned} \quad (2)$$

where “ $*$ ” indicates convolution. By hypothesis, $H(z)$ and $G(z)$ are approximate inverses of one another within a delay of D , so that (2) becomes

$$y[n] \approx x[n - D] - g[n] * e[n] \quad (3)$$

and the error $r[n]$ for the overall system can be approximated as

$$r[n] \approx g[n] * e[n]. \quad (4)$$

By assumption, $x[n]$ is stationary, so that $e[n]$ and $r[n]$, which are derived as time-invariant (but nonlinear) functions of $x[n]$, are likewise stationary [12]. Let the power spectra of $e[n]$ and $r[n]$ be denoted $S_{ee}(\omega)$ and $S_{rr}(\omega)$, respectively, where ω is radian frequency. In terms of these power spectra, (4) can be rewritten [12]

$$S_{rr}(\omega) \approx |G(e^{j\omega})|^2 S_{ee}(\omega); \quad -\pi \leq \omega < \pi. \quad (5)$$

It is assumed that both prefilter and postfilter are approximately all-pass, meaning that their magnitude-frequency responses are nearly flat over the full spectrum. Consistent with this, it can be assumed without loss of generality that $|G(e^{j\omega})| = 1$ for $-\pi \leq \omega < \pi$ (whatever scaling factor is needed to make this true can be canceled by an appropriate gain in the prefilter). Thus, (5) reduces to

$$S_{rr}(\omega) \approx S_{ee}(\omega), \quad (6)$$

which implies that the mean-square value of $r[n]$ is nearly equal to that of $e[n]$. That is, the system’s MSE is nearly equal to that incurred by quantizing the intermediate (prefiltered) signal $u[n]$.

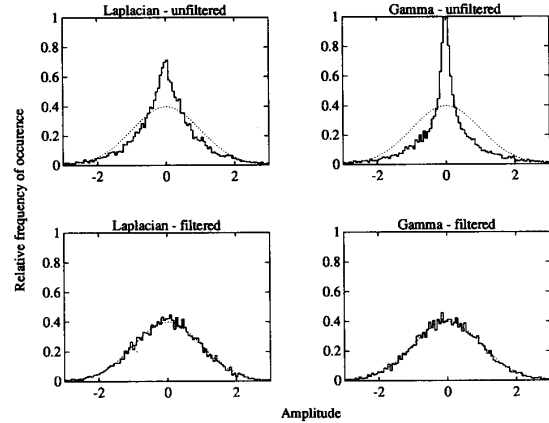


Fig. 3. Histograms of original and filtered sources, based on 10 000 pseudo-random samples. A 31 tap filter was used in the Laplacian case, while a 60 tap filter was used in the gamma case. For reference, an appropriately scaled version of the Gaussian PDF is also shown (dotted curve).

III. STATISTICAL CHARACTERIZATION OF THE INTERMEDIATE SIGNAL

Since each sample in $u[n]$ is the (weighted) sum of *independent* random variables, its probability density function is approximately Gaussian, provided that the sum includes a sufficient number of variables, each with nonnegligible but not disproportionately large weight (by Liapounov’s central limit theorem [13]). These conditions will be satisfied when the impulse response of the prefilter $h[n]$ has significantly nonzero values distributed over a sufficiently long interval. It is not difficult to show that for a memoryless input, a sufficient condition for a sequence of FIR filters,

$$P_N(z) = \sum_{i=1}^N \rho_{N,i} z^{-1}$$

to asymptotically produce marginally Gaussian output is that the sequence of numbers $\{\min_i |\rho_{N,i}|\}$ be bounded and also bounded away from zero. In this paper, a filter that has this property will be called time-dispersive.

The most convenient measure of the extent to which the pdf of the source is modified is the observed performance of simple quantization of the intermediate signal under the assumption of a Gaussian distribution. It has been found experimentally that for a Laplacian source, FIR time-dispersive prefilters and postfilters of length 30 are sufficiently long to yield a signal-to-quantization MSE ratio within 0.2 dB of the best possible (that for a Gaussian source) at rates up to 5 bits per sample. In the case of a gamma source, a length of 120 is required for the same level of performance. However, even when much shorter filters are used, a significant improvement over direct quantization results (see Section V). Design of appropriate prefilters and postfilters of a given length is discussed in Section IV.

Another way to gauge the extent to which the pdf of the input is modified is to examine histograms. Fig. 3 shows histograms based on 10 000 samples from simulated sources, before and after prefiltering.

A legitimate objection to the foregoing analysis is that in many applications, the assumption of independence of successive source samples is unjustified, so that prefiltering may not make the pdf approximately Gaussian. In fact, it is easy to construct a source for which prefiltering makes the distribution appear *less* Gaussian—for example, simply filter a gamma source by $G(z)$, and use the result as $x[n]$. In all such cases, the previous analysis can be made to apply if the source samples are rearranged or scrambled in a pseudorandom manner prior to prefiltering, and subsequently restored to their original ordering after postfiltering. Note that scrambling and inverse scrambling are linear and energy-preserving operations, so that the analysis of Section II applies, and quantization MSE is preserved. However, scrambling and inverse scrambling necessarily introduce considerable delay, and therefore may not be appropriate in some applications.

IV. DESIGN OF TIME-DISPERSIVE PREFILTERS AND POSTFILTERS

It is desired that the prefilter and postfilter be approximate inverses of one another, that their impulse responses have envelopes that extend sufficiently over time, and that their magnitude-frequency responses be approximately flat. One approach to obtaining such filters is to begin with “initial guess” filters that have *roughly* the right properties, then refine these by numerical optimization. In particular, the following procedure has proven to be successful.

Begin with a windowed “chirp” signal (swept sinusoid) as an initial guess for the impulse response of the prefilter, and use the same chirp, but time-reversed, as the initial guess for the impulse response of the postfilter. To simplify notation in the present section, the postfilter is not required to be casual, and the delay D of the cascade is taken to be zero. Each of the two initial guess filters has the desired property that the energy in its impulse response is distributed over the entire region of support; that is, the filters are time-dispersive. In order to ensure that the remaining requirements are met—that the prefilter and postfilter be approximate inverses of each other, and that each have an approximately flat magnitude-frequency response—a numerical procedure is used to modify the initial guess filters to minimize the total square difference between the convolution of $h[n] * g[n]$ and the unit-sample sequence $\delta[n]$. That is, a local minimum is sought of the objective function

$$E = \sum_{n=-\infty}^{\infty} \left[\sum_{k=-\infty}^{\infty} h[k]g[n-k] - \delta[n] \right]^2 \quad (7)$$

over the joint space of prefilter and postfilter coefficients $\{h[n], g[n]\}$, beginning the search at the specified initial guess.

Observe that the initial guess filters, being time-reversed versions of each other, have identical magnitude-frequency responses. By maintaining this relationship throughout the optimization, so that the optimized filters also end up as time-reversed versions of each other, the magnitude-frequency response of each optimized filter can be made to be approximately flat. This follows because the magnitude-frequency

response of the cascade—which is the product of the individual responses—must be flat if the two filters are to be inverses of each other. The time-reversed relationship can either be maintained explicitly by adding a simple constraint (i.e., optimizing over only one of the filters and fixing the other according to the time-reversed relationship), or else the symmetry of the objective function can be relied upon to maintain the relationship from the initial guess. The latter approach was found to work consistently in the present investigation.

It is natural to question the existence of local minima, convergence issues, and so on, however, such a formal treatment of the optimization problem is avoided here, on the grounds that in practice, a local minimum seems to be obtainable quickly and consistently using any of a variety of well-known optimization procedures.

Figs. 4 and 5 show the characteristics of the prefilter and postfilter in a 60 tap design example, before and after optimization, respectively. Also shown in each figure is the convolution of $h[n]$ and $g[n]$ and the corresponding magnitude-frequency response of the cascade of the two filters. In this example, the chirp initial guess filters were taken to be

$$h[n] = g[-n] = A \sin(\pi(n^2 - n)/120), \quad n = 0, \dots, 59$$

where A was chosen to make the sum of the squares of the coefficients in each filter unity. Note that the optimized filters (Fig. 5) have the desired properties: the prefilter and postfilter are approximate inverses of each other, have nearly flat magnitude-frequency responses, and have impulse responses with significant energy distributed well over the entire region of support.

V. EXPERIMENTAL RESULTS

The dependence of the performance of the proposed system on the length of the filters is illustrated in Fig. 6, for the range of 5–37 taps. To obtain each point in the graphs, the sources were simulated using techniques described by Knuth [14] and the performance of the proposed quantization system was measured for 10000 samples. Observe that for both sources, even very short filters yield a considerable improvement in performance. Although not shown in the figure, it was found that as the filters are made longer than 37 taps, the improvement in performance continues to be noticeable for the gamma source, but not for the Laplacian source.

In obtaining the remaining experimental results presented in this section, filters of length 31 and 60 were used in the Laplacian and gamma cases, respectively. Also, unless otherwise stated, all measurements were based on 10000 samples.

Fig. 7 shows the SNR in dB as a function of the number of quantization bits for simulated Laplacian and Gamma sources, for both simple quantization and the filter-based quantization scheme. Also shown are samples of the rate-distortion functions for these two sources. Note that the performance is, as expected, approximately that of simple quantization of a Gaussian source. The improvement over quantization without filtering is particularly significant at low bit-rates; in fact, by comparing the results with those presented in [2] it

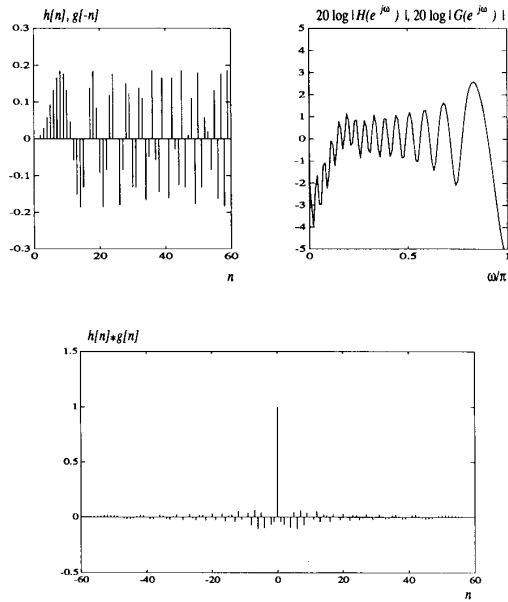


Fig. 4. Prefilter and postfilter characteristics *before* optimization (see Section IV).

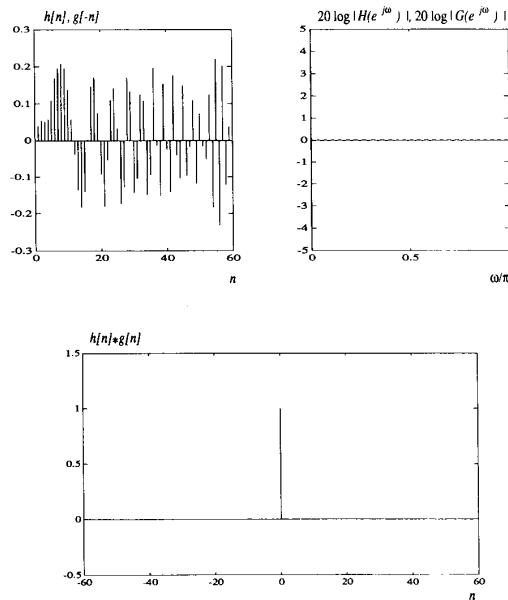


Fig. 5. Prefilter and postfilter characteristics *after* optimization (see Section IV).

can be concluded that at 1 bit/sample, the improvement over direct quantization obtained by prefiltering and postfiltering is roughly the same as would be obtained by three-dimensional vector quantization.

It should be noted that although one advantage of the proposed technique over simple quantization is the reduced MSE for some sources, alternative existing schemes such as vector quantization can provide even better performance. The

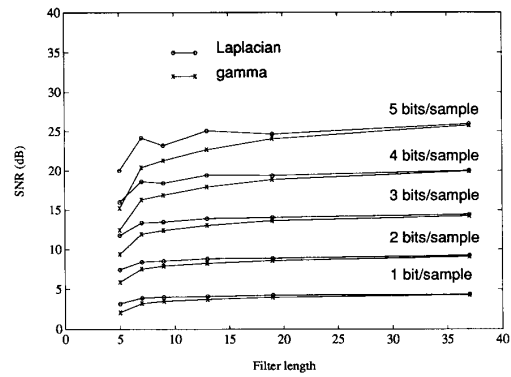


Fig. 6. Experimentally determined dependence of performance of the proposed system on the length of the filters. The procedure described in Section IV was used to design the filters.

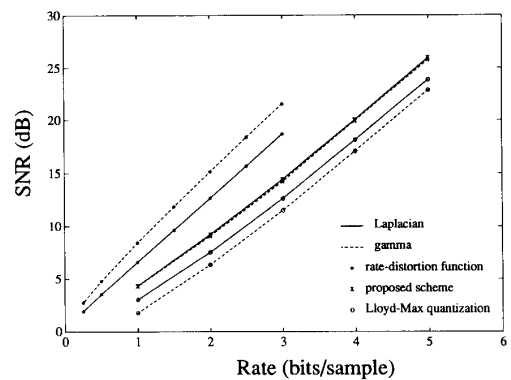


Fig. 7. Experimental performance (in terms of mean-square error) of the proposed scheme relative to direct Lloyd-Max quantization for a Laplacian and a gamma source. Also shown are samples of the corresponding rate-distortion function, computed via the Blahut algorithm [16].

more important feature, however, is the proposed scheme's robustness to errors in modeling the input. Specifically, the prefiltering operation will tend to make *any* memoryless source appear more Gaussian, so that the performance of the system does not depend critically on accurate modeling of the input pdf. Fig. 8 illustrates the relative insensitivity of the proposed technique to modeling errors. Shown is the performance of the proposed system and that of direct quantization, when each source is mistakenly modeled as the other. Note that the performance of direct quantization is reduced because of the mismatch, while that of the proposed system is unaffected. This implies that the proposed system can be used with some confidence even when relatively little is known about pdf of the source. In some applications, this robustness is more significant than reduction of mean-square error. The most notable prior work in robust simple quantization is by Bath and Vandelinde [15]. In their approach, a minimum level of MSE performance is guaranteed so long as the input pdf belongs to a certain class; however, that performance is considerably worse than the performance of Lloyd-Max quantization of a Gaussian source. In contrast, for the same class of input pdf, the performance of the quantization scheme proposed

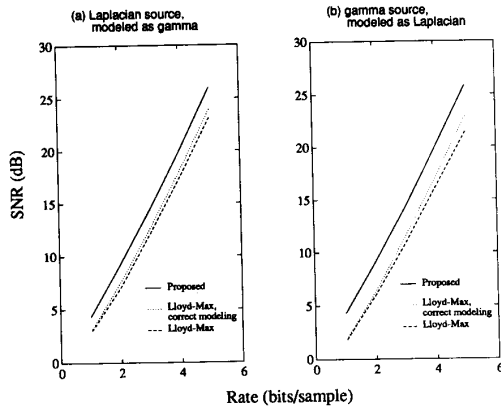


Fig. 8. Experimental mean-square error performance of the proposed system and that of direct quantization in the case of quantizer mismatch. For reference, also shown is the mean-square error of direct quantization when the input is correctly modeled (dotted line). In (a) the source is Laplacian distributed, but modeled as gamma. In (b) the source is gamma distributed, but modeled as Laplacian.

in this paper is always *roughly equal* to that of Lloyd–Max quantization of a Gaussian source.

Quantizers often play the role of fundamental building blocks in more complex source coding systems. In such systems it is often assumed that the input to the quantizer is Gaussian and memoryless. In such systems, the scheme presented in this paper may serve a useful role for the quantization component of the overall system. For example, in [8] the authors point out that trellis coded quantization (TCQ) performs better for the Gaussian memoryless source than any other known coding techniques. By embedding the dispersive filtering scheme inside TCQ, it may be possible to provide a high performance quantization system that is very robust to changing source statistics. A similar application might be found in DPCM systems designed specifically for Gauss-Markov sources, in which a linear prediction filter is assumed to produce a Gaussian i.i.d. output.

Finally, it should be noted that in evaluating the proposed system's performance, there are tradeoffs between delay, complexity, and robustness. For example, the scheme proposed in this paper requires two FIR filterings, each requiring a number of arithmetic multiplies and additions in proportion to the filter lengths D . In contrast, the Bath and Vandelinde quantizer is

less complex to implement but does not provide the same degree of robustness. Hence, the robustness of the scheme presented in the present paper is achieved at the cost of some extra complexity and delay.

VI. CONCLUSION

A novel technique for quantization of memoryless sources has been presented. Experimental results have confirmed that the system performs as expected. The technique results in a reduction in mean-square quantization error for certain sources and offers relative insensitivity to errors in modeling the input distribution.

REFERENCES

- [1] M. D. Paez and T. H. Glisson, "Minimum mean squared error quantization in speech PCM and DPCM systems," *IEEE Trans. Commun.*, vol. COM-20, pp. 225–230, Apr. 1972.
- [2] T. R. Fischer and R. M. Dicharry, "Vector quantizer design for memoryless Gaussian, gamma, and Laplacian sources," *IEEE Trans. Commun.*, vol. COM-32, pp. 1065–1069, Sept. 1984.
- [3] N. Jayant and P. Noll, *Digital Coding of Waveforms, Principles and Applications to Speech and Audio*. Englewood Cliffs, NJ: Prentice Hall, 1984.
- [4] L. R. Rabiner and R. Schafer, *Digital Coding of Speech Signals*. Englewood Cliffs, NJ: Prentice-Hall, 1978.
- [5] S. P. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 129–137, Mar. 1982.
- [6] J. Max, "Quantizing for minimum distortion," *IRE Trans. Inform. Theory*, pp. 7–12, Mar. 1960.
- [7] T. Berger, "Minimum entropy quantizers and permutation codes," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 759–765, Mar. 1982.
- [8] M. W. Marcellin and T. R. Fischer, "Trellis coded quantization of memoryless and Gauss–Markov sources," *IEEE Trans. Commun.*, vol. 38, pp. 82–93, Jan. 1990.
- [9] A. J. Viterbi and J. K. Omura, "Trellis encoding of memoryless discrete-time sources with a fidelity criterion," *IEEE Trans. Info. Theory*, vol. IT-20, pp. 325–332, May 1974.
- [10] H. W. Strube, "How to make an all-pass filter with a desired impulse response," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-30, pp. 336–337, Apr. 1982.
- [11] Zenith Electronics Corporation, *Zenith spectrum compatible HDTV system*. Glenview, IL, 1988.
- [12] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*. New York: McGraw-Hill, 1984.
- [13] J. B. Thomas, *An Introduction to Applied Probability and Random Processes*. New York: Wiley, 1971.
- [14] D. E. Knuth, *The Art of Computer Programming*, 3. Reading, MA: Addison-Wesley, MA: 1973.
- [15] W. G. Bath and V. D. Vandelinde, "Robust memoryless quantization for minimum signal distortion," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 296–297, Mar. 1982.
- [16] R. E. Blahut, "Computation of channel capacity and rate-distortion functions," *IEEE Trans. Inform. Theory*, vol. IT-18, pp. 460–473, July 1972.