# *Special Issue on ITG Conference on Source and Channel Coding*

# Performance of multistage vector quantization in hierarchical coding

Christoph Erdmann* and Peter Vary

*Institute of Communication Systems and Data Processing, Aachen University of Technology, Templergraben 55, D-52056 Aachen, Germany*

## SUMMARY

How much performance penalty does a hierarchical coder based on multistage vector quantization (MSVQ) suffer compared to a non-hierarchical coder based on fixed rate VQ? In this paper, the above question shall be answered from a rate-distortion theoretic perspective. We analyze several results from high-rate or asymptotic quantization theory and use them to specify an upper bound on the MSVQ penalty in terms of mean square error (MSE) distortion. The theoretical results are used to gain more analytic insight in hierarchical coding systems based on a multistage coding approach. Although entirely based on high-rate assumptions, in practice this bound also applies for relatively small rates as shown by experiment. Copyright © 2004 AEI.

## 1. INTRODUCTION

Usually, hierarchical coding based on multistage vector quantization (MSVQ) involves a performance penalty when compared to non-hierarchical coding based on fixed rate VQ. This statement has some intuitive appeal. Also it is well-known from rate-distortion theory [1] that the performance bound for quantization of a random variable can be approached by using VQ of arbitrary large dimension or block length.

It is sometimes argued that the results provided by rate-distortion theory are of limited relevance for practical quantization tasks since both search complexity and memory requirements grow exponentially with the block length. A contrasting approach is to fix the block length and assume the bit rate (and hence the codebook size) to be large, which is the subject of high-rate or asymptotic quantization theory [2]. High-rate theory thus provides results on VQ performance for any block length which actually show that even in case of finite block lengths VQ gets closer to the rate-distortion bound for a given source than any other coding scheme that breaks down the quantization task into several sub-tasks of reduced block length.

To develop a basic understanding how different quantization schemes utilize such characteristic signal properties as dimension, correlation between samples and the probability density function (PDF) for efficient quantization in a rate-distortion sense, we review some results from rate-distortion theory and high-rate theory. Although the majority of these *theoretical tools* are well-known from the literature and have been frequently used to describe the performance of fixed rate VQ, in this paper they are used for the first time to describe the performance respectively the performance penalty of MSVQ.

Our particular treatment of MSVQ is motivated by the fact that it represents a typical example of so-called *sequential search product code* technique [3], which is frequently used in practical implementations of hierarchical coding systems.

In this paper, we address the performance penalty of MSVQ compared to fixed rate VQ using a mean square error distortion criterion. In Section 2, we review some important results from rate-distortion and high-rate theory that are used to describe the theoretical performance of VQ. In Section 3, the concept of MSVQ is introduced prior to the derivation of an upper bound on the MSVQ penalty. To see how this bound applies in practice for relatively

---

* Correspondence to: Christoph Erdmann, Institute of Communication Systems and Data Processing, Aachen University of Technology, Templergraben 55, D-52056 Aachen, Germany. E-mail: erdmann@ind.rwth-aachen.de

small rates, two examples of MSVQ are discussed in Section 4.

## 2. THEORETICAL VQ PERFORMANCE

Rate-distortion theory defines optimum quantizer performance for a given source and mean square error (MSE) distortion by a distortion-rate function (DRF) $\mathfrak{D}(R)$, that describes the minimum bit rate $R$ which is required for quantization with a given MSE distortion $\mathfrak{D}$. One of the basic statements resulting from rate-distortion theory is that by using a vector quantizer one can in principle approach the DRF of any given source arbitrarily closely by increasing the vector dimension $d$. Thus, $\mathfrak{D}(R)$ is not merely a lower bound on the achievable quantizer distortion, it is actually achievable by VQ with high dimension.

The DRF $\mathfrak{D}(R)$ has two important properties: (1) it is monotone decreasing with $R$ and (2) it is convex. Furthermore, for the MSE distortion in decibels, $\mathfrak{D}(R)$ decreases at a rate of 6.02 dB/bit for large $R$. For a zero-mean, memoryless Gaussian source with variance $\sigma^2$, the DRF with MSE distortion normalized to $\sigma^2$ is known as

$$\mathfrak{D}_G(R) = 10 \log \frac{\mathfrak{D}_G(R)}{\sigma^2} = -6.02\,R \quad [\text{dB}] \qquad (1)$$

where $\mathfrak{D}$ denotes the normalized MSE distortion in decibels, which shall be used to describe quantizer performance for the remainder of this paper.

There are only scant explicit $\mathfrak{D}(R)$ results for sources other than the memoryless Gaussian source. Yet, lower and upper bounds exist so that

$$\mathfrak{D}_{SLB}(R) \leqslant D(R) \leqslant \mathfrak{D}_G(R) \qquad (2)$$

whereas $\mathfrak{D}(R)$ is known to be upper bound by the DRF $\mathfrak{D}_G(R)$ for a Gaussian source. The lower bound $\mathfrak{D}_{SLB}(R)$ is the *Shannon lower bound* (SLB) [1] which for MSE distortion is given by

$$\mathfrak{D}_{SLB}(R) = \frac{1}{2\pi e} 2^{2(h(x)-R)} \qquad (3)$$

where

$$h(x) = -\int_{-\infty}^{\infty} p_x \log_2 p_x \qquad (4)$$

is the *differential entropy* of the memoryless source with PDF $p_x$. According to the special role that Gaussian sources play in bounding the performance of coding systems, also $h(x)$, for $x$ being a random variable with

variance $\sigma^2$, is upper bound by the differential entropy $h_G(x)$ of a Gaussian PDF with equal variance $\sigma^2$ (e.g. Reference [4]), thus,

$$h(x) \leqslant h_G(x) = \frac{1}{2} \log_2 \left( 2\pi e \sigma^2 \right) \qquad (5)$$

For many sources the SLB is achievable only as $R \to \infty$. With Equations (3) and (5), we can write $\mathfrak{D}_{SLB}(R)$ as a normalized MSE distortion in decibels

$$\mathfrak{D}_{SLB}(R) = -6.02R - 6.02(h_G(x) - h(x))\,[\text{dB}] \qquad (6)$$

Hence, with Equation (1) the asymptotic difference in MSE distortion between the Gaussian upper bound and the SLB amounts to

$$\begin{aligned} \Delta\mathfrak{D}_{G,SLB} &= \mathfrak{D}_G(R) - \mathfrak{D}_{SLB}(R) \\ &= 6.02(h_G(x) - h(x)) \quad [\text{dB}] \end{aligned} \qquad (7)$$

Since $h_{G(x)} > h(x)$, the SLB $\mathfrak{D}_{SLB}(R)$ is lower than the Gaussian DRF by an amount equal to the difference between the differential entropies $h_G(x)$ and $h(x)$ (in bits) multiplied by 6.02 dB/bit. Equation (6) clearly indicates that the asymptotic behavior of the DRF for many PDFs is expected to decrease at a rate of 6.02 dB/bit as $R \to \infty$.

Table 1 shows the asymptotic difference $\Delta\mathfrak{D}_{G,SLB}$ between the Gaussian upper bound and the SLB according to Equation (7) for four PDFs that are common models used for certain signal distributions.

The plots in Figure 1 illustrate the DRF $\mathfrak{D}_\Gamma(R)$ for a Gamma PDF along with the Gaussian upper bound $\mathfrak{D}_G(R)$ and the SLB $\mathfrak{D}_{SLB}(R)$. Noll and Zelinski [5] obtained the DRF $\mathfrak{D}_\Gamma(R)$ by numerical calculations using *Blahut's algorithm* [6]. The Gamma PDF was chosen among the four PDFs from Table 1 because it shows very clearly the departure of $\mathfrak{D}(R)$ from the Gaussian DRF.

Table 1. Four common PDFs and their asymptotic difference in MSE distortion between the Gaussian DRF $\mathfrak{D}_G(R)$ and the Shannon lower bound $\mathfrak{D}_{SLB}(R)$ in decibels.

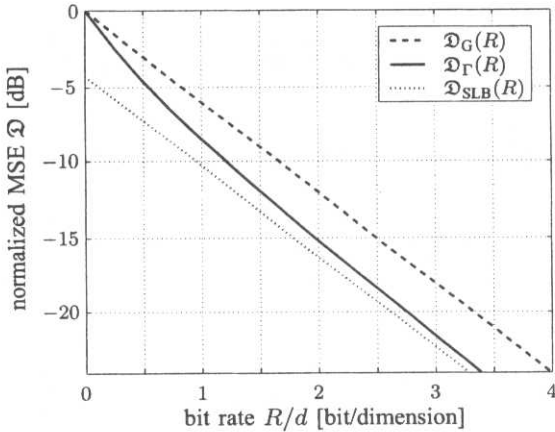| PDF | $p_x$ | $\Delta\mathfrak{D}_{G,SLB}[\text{dB}]$ |
|---|---|---|
| Gaussian (G) | $\frac{1}{\sqrt{2\pi\sigma^2}} \exp[-x^2/2\sigma^2]$ | 0 |
| Uniform (U) | $\frac{1}{2\sqrt{3\sigma^2}}, \quad \|x\| \leqslant \sqrt{3\sigma^2}$ | 1.53 |
|  | $0, \quad \text{otherwise}$ | |
| Laplacian (L) | $\frac{1}{\sqrt{2\sigma^2}} \exp\left[-\sqrt{2}\|x\|/\sigma\right]$ | 0.63 |
| Gamma ($\Gamma$) | $\frac{\sqrt[4]{3}}{\sqrt{8\sigma\pi\|x\|}} \exp\left[-\sqrt{3}\|x\|/2\sigma\right]$ | 4.27 |

Figure 1. Distortion-rate function $\mathfrak{D}_\Gamma(R)$ (in decibels, MSE distortion) for a memoryless Gamma source. $\mathfrak{D}_\Gamma(R)$ is upper bounded by the Gaussian DRF $\mathfrak{D}_G(R)$ and lower bounded by the SLB $\mathfrak{D}_{SLB}(R)$ (from Reference [5]).

Note that the $\mathfrak{D}_\Gamma(R)$ curve in Figure 1 is monotone decreasing and convex. Also, as $R$ increases beyond a few bits, $\mathfrak{D}_\Gamma(R)$ decreases at a rate of 6.02 dB/bit which is the slope of the upper and lower bounds.

## 2.1. Scalar quantization results

From a quantizer design perspective, a scalar quantizer (SQ) such as the well-known *Lloyd–Max* quantizer (LMQ) [7, 8] can be viewed as a special case of one-dimensional VQ. Let $R = \log_2 N$ be the bit rate given in bit/sample for SQ with $N$ reconstruction levels. For large $R$, an asymptotic formula for the normalized MSE distortion of LMQ is provided in Reference [9] that can be written in decibels as

$$\mathfrak{D}_{LMQ}(R) = -6.02R + \mathfrak{F}_{LMQ}(p_x) \quad [dB] \qquad (8)$$

where $\mathfrak{F}_{LMQ}(p_x)$ is a constant that depends on the PDF shape. Note again the $-6.02$ dB/bit behavior for large $R$. However, as the first term in Equation (8) equals $\mathfrak{D}_G(R)$ from Equation (1), we can modify Equation (8) to obtain

$$
\begin{aligned}
\mathfrak{F}_{LMQ}(p_x) &= \mathfrak{D}_{LMQ}(R) - \mathfrak{D}_G(R) \\
&= 10 \log \epsilon^2(p_x) \quad [dB]
\end{aligned}
\qquad (9)
$$

Hence, the quantity $\epsilon^2(p_x)$, or $\mathfrak{F}_{LMQ}(p_x)$ in decibels, describes the PDF shape dependent performance of practical PDF-optimized scalar quantizers such as LMQ in terms

Table 2. Performance of LMQ for four common PDFs.

| PDF | $\mathfrak{D}_{LMQ}(R) - \mathfrak{D}_{SLB}(R)$ [dB] | $\mathfrak{F}_{LMQ}(p_x)$ [dB] | $\epsilon^2(p_x)$ |
|---|---|---|---|
| G | 4.347 | 4.347 | 2.721 |
| U | 1.53 | 0 | 1 |
| L | 7.17 | 6.537 | 4.505 |
| $\Gamma$ | 11.82 | 7.547 | 5.685 |

of normalized MSE distortion relative to the Gaussian DRF $\mathfrak{D}_G(R)$. The parameter $\epsilon^2(p_x)$ is also known as *SQ performance factor* [10].

Table 2 shows the asymptotic difference between the normalized MSE distortion of LMQ $\mathfrak{D}_{LMQ}(R)$ and the Shannon lower bound $\mathfrak{D}_{SLB}(R)$. Note that for high rates $R \to \infty$ the DRF $\mathfrak{D}(R)$ approaches $\mathfrak{D}_{SLB}(R)$. The difference $\mathfrak{D}_{LMQ}(R) - \mathfrak{D}_{SLB}(R)$ represents the maximum performance that VQ with large dimension $d \to \infty$ can potentially gain over LMQ. $\mathfrak{D}_{LMQ}(R) - \mathfrak{D}_{SLB}(R)$ is thus called the asymptotic *VQ gain*.

Columns 3 and 4 of Table 2 list the PDF dependent performance of LMQ relative to the DRF $D_G(R)$ for the Gaussian PDF, expressed either by the quantity $\mathfrak{F}_{LMQ}(p_x)$ in decibels or the corresponding SQ performance factor $\epsilon^2(p_x)$ on a linear scale.

Having the most peaked shape among the four model PDFs, the Gamma PDF results in the highest MSE distortion when using LMQ, although it has the lowest DRF $\mathfrak{D}_\Gamma(R)$. This result is also shown by the various $\mathfrak{D}_{LMQ}(R)$-curves in Figure 2, which were obtained from particularly designed LMQs. Therefore each scalar quantizer has been optimized to one of the four PDFs by using
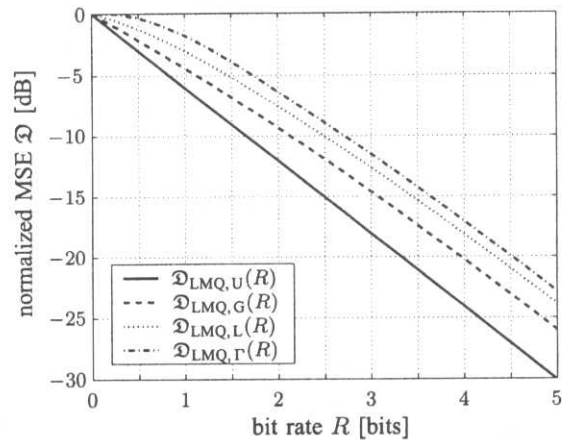


Figure 2. Normalized MSE distortion $\mathfrak{D}_{LMQ}(R)$ for a uniform, Gaussian, Laplacian and Gamma distributed memoryless sources using LMQ with $R = \log_2 N$ bit/sample.

the LBG algorithm [11], which in the scalar case yields the same quantizer as Lloyd's algorithm for PDF-optimized quantizer design [7]. (For comparison see also $\mathfrak{D}_{LMQ}(R)$-values in Table 4.4 in Reference [10]).

## 2.2. Asymptotic VQ results

The above given asymptotic formulas for SQ performance have analogues for VQ. In References [12, 13] it was shown that under the assumption of high rates $R$ given in bit/vector a $d$-dimensional VQ with $N = 2^R$ reproduction vectors attains a minimum MSE distortion of

$$\mathfrak{D}_{VQ}(d, N) = C_Q(d) N^{-\frac{2}{d}} \|p_{\mathbf{x}}\|_{d/(d+2)} \quad (10)$$

where $C_Q(d)$ is the *coefficient of quantization* that describes how well cells can be packed in $\mathbb{R}^d$ independent of the PDF [14]. $\|p_{\mathbf{x}}\|_r$ denotes the $r$th norm of the d-dimensional PDF $p_{\mathbf{x}}$ according to

$$\|p_{\mathbf{x}}\|_r = \left[ \iint_{\mathbb{R}^d} (p_{\mathbf{x}})^r \, d\mathbf{x} \right]^{\frac{1}{r}} \quad (11)$$

For $N$-level VQ with $N = 2^R$, Equation (10) may by written in decibels as

$$\mathfrak{D}_{VQ}(d, R) = -\frac{6.02}{d} R + \mathfrak{F}_{VQ}(d, p_{\mathbf{x}}) \quad [dB] \quad (12)$$

Analog to the term $\mathfrak{F}_{LMQ}(p_x)$ in Equation (8), the term $\mathfrak{F}_{VQ}(d, p_{\mathbf{x}})$ describes the performance of PDF-optimized VQ relative to DRF for the Gaussian PDF, but now depending on the joint PDF $p_{\mathbf{x}}$ and the dimension $d$, i.e.

$$\mathfrak{F}_{VQ}(d, p_{\mathbf{x}}) = 10 \log \left[ C_Q(d) \|p_{\mathbf{x}}\|_{d/(d+2)} \right] \quad (13)$$

Comparing this result with the scalar case from Equation (9), we can readily see that $\mathfrak{F}_{VQ}(d, p_{\mathbf{x}})$ is not merely a generalization of $\mathfrak{F}_{LMQ}(p_x)$. Beyond the dependence on the PDF shape, the step from one dimension to multiple dimensions $d$ causes the quantity $\mathfrak{F}_{VQ}(d, p_{\mathbf{x}})$ to include further signal properties that may have an impact on the VQ performance such as vector dimension and linear as well as non-linear dependencies between vector components.

## 3. MULTISTAGE VQ

Multistage VQ, which is sometimes also called *cascaded* VQ, is quite common for quantization of, for example, LP parameters in predictive speech coding. According to the concept of successive refinement, the basic idea of MSVQ
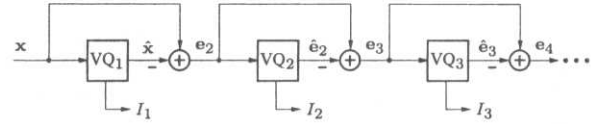


Figure 3. Multistage vector quantization (MSVQ).

as shown in Figure 3 is to divide encoding with $N$ reproduction levels into $\nu$ successive stages with $N_i$ reproduction levels each, so that $N = \prod_{i=1}^{\nu} N_i$.

The first stage VQ$_1$ performs a relatively crude quantization of the input vector $\mathbf{x}$ using a small codebook of size $N_1$. Then, a second stage quantizer VQ$_2$ with $N_2$ reproduction levels operates on the error vector $\mathbf{e}_2$ between the original $\mathbf{x}$ and quantized reproduction $\hat{\mathbf{x}}$ of the first stage. The quantized error $\hat{\mathbf{e}}_2$ then provides a second approximation to the original input vector thereby leading to a refined representation of the input. A third stage quantizer VQ$_3$ with $N_3$ reproduction levels may then be used to quantize the second stage error vector $\mathbf{e}_3$ to provide a further refinement and so on.

From a hierarchical coding perspective, each coding stage consists of pure VQ with $N_i$ reproduction levels, where the input vector for the next coding stage is the corresponding quantization error. It is important to note that all stages operate on the full dimensionality $d$. MSVQ thus realizes only a decomposition in terms of reduced codebook sizes but not in terms of reduced dimensionality.

### 3.1. Performance of MSVQ

To analyze the MSE distortion of MSVQ, let us consider two quantizers VQ$_1$ and VQ$_\Sigma$ of equal dimension $d$ but with different bit rates $R_1$ and $R_\Sigma$, where $R_1 < R_\Sigma$. The bit rates are assumed to be sufficiently large, so that the individual MSE distortions are given by the asymptotic formula from Equation (12). Both, VQ$_1$ and VQ$_\Sigma$ are separately applied to the same vector $\mathbf{x}$ with a given but arbitrary PDF $p_{\mathbf{x}}$. As $\mathfrak{F}_{VQ}(d, p_{\mathbf{x}})$ depends only on the properties of the input vector $\mathbf{x}$, while being independent of the quantizer, the respective value of $\mathfrak{F}_{VQ}(d, p_{\mathbf{x}})$ is equal for both VQs, so that the difference of the individual MSE distortions in decibels is given by

$$\begin{aligned} \mathfrak{D}_{VQ}(d, R_\Sigma) - \mathfrak{D}_{VQ}(d, R_1) &= -\frac{6.02}{d}(R_\Sigma - R_1) \\ &= -\frac{6.02}{d} R_\Delta \quad [dB] \end{aligned} \quad (14)$$

This result may seem trivial, but it gives rise to an interesting interpretation when we consider MSVQ as shown in

Figure 3 with only two stages, with quantizer $VQ_1$ in the first stage, and the quantizer $VQ_2$ in the second stage employing a bit rate of $R_\Delta$ bit/vector.

Certainly, the overall MSE distortion of such two-stage MSVQ can never fall below the MSE distortion of $VQ_\Sigma$, i.e.

$$\mathfrak{D}_{VQ}(d, R_\Sigma) \leqslant \mathfrak{D}_{VQ}(d, R_1) + \mathfrak{D}_{VQ}(d, R_\Delta) \, [\text{dB}] \quad (15)$$

Otherwise, this would imply that the same performance can be exceeded by sharing a fixed quota of bits over several VQs, than by using all the bits in one single $VQ_\Sigma$ while at the same time achieving a great complexity reduction. With Equation (14) we obtain the following condition for the MSE distortion $\mathfrak{D}_{VQ}(d, R_\Delta)$ of $VQ_2$:

$$-\frac{6.02}{d} R_\Delta \leqslant \mathfrak{D}_{VQ}(d, R_\Delta) \quad [\text{dB}] \quad (16)$$

Hence, we can conclude that the MSE distortion $\mathfrak{D}_{VQ}(d, R_\Delta)$ of $VQ_2$ is bounded below by the DRF $\mathfrak{D}_G(R)$ for a Gaussian PDF, even if $R_\Delta$ is sufficiently large so that $VQ_2$ itself falls into the high-rate case.

(1) *Memoryless Sources*: Let us firstly consider the $d$-dimensional input vector $\mathbf{x}$ to be memoryless, so that its joint PDF $p_\mathbf{x}$ is completely specified by its respective marginal PDF $p_x$ according to $p_\mathbf{x} = (p_x)^d$. In this case, we have $\mathfrak{F}_{VQ}(d, p_\mathbf{x}) \equiv \mathfrak{F}_{VQ}(d, p_x)$. The distortion of $VQ_2$ is thus entirely determined by the shape of the marginal PDF of the quantization error vector $\mathbf{e}_2$, (i.e., the output of the first stage) and the vector dimension $d$. Recall that in MSVQ the vector dimension remains the same for all stages.

Considering for example the same Gamma distributed source as in Figure 1, Figure 4 illustrates the asymptotic VQ gain $\mathfrak{D}_{LMQ}(R) - \mathfrak{D}_{SLB}(R)$ from Table 2 that describes the potential improvement in MSE distortion by VQ

relative to LMQ, for high bit rates $R$ and large dimensions $d$. On the other hand the term $\mathfrak{F}_{LMQ}(p_x)$ specifies the MSE distortion of LMQ relative to the Gaussian DRF (see Equation (9) and Table 2). Hence, as $d \rightarrow \infty$ the limiting value of the MSE distortion of VQ relative to the Gaussian DRF is given by

$$\lim_{d \rightarrow \infty} \mathfrak{F}_{VQ}(d, p_x) = \mathfrak{F}_{LMQ}(p_x)$$
$$- (\mathfrak{D}_{LMQ}(R) - \mathfrak{D}_{SLB}(R)) \, [\text{dB}] \quad (17)$$

and with the expression for $\mathfrak{F}_{LMQ}(p_x)$ from Equation (9):

$$\lim_{d \rightarrow \infty} \mathfrak{F}_{VQ}(d, p_x) = \mathfrak{D}_{SLB}(R) - \mathfrak{D}_G(R)$$
$$= -\Delta \mathfrak{D}_{G,SLB} \quad [\text{dB}] \quad (18)$$

The limiting value of $\mathfrak{F}_{VQ}(d, p_x)$ is thus equal to the negative asymptotic difference $-\Delta \mathfrak{D}_{G,SLB}$ between the Gaussian DRF and the Shannon lower bound from Equation (7). These implications are also illustrated in Figure 4.

For finite dimensions $d \geqslant 1$, Figure 5 shows four distinct plots of $\mathfrak{F}_{VQ}(d, p_x)$ corresponding to the four model PDFs from Table 1 [15]. Several observations can be made on the $\mathfrak{F}_{VQ}(d, p_x)$-curves. First, we note that for large dimensions, and in accordance to Equation (18), the limiting values of the $\mathfrak{F}_{VQ}(d, p_x)$-curves are indeed the negative values of $\Delta \mathfrak{D}_{G,SLB}$ as listed in Table 1.

A second observation relates to the case that the quantization error vector $\mathbf{e}_2$ with marginal PDF $p_{e_2}$ is to be quantized by $VQ_2$. According to Equation (16) the Gaussian DRF represents the minimum achievable MSE distortion of $VQ_2$. Substitution of $\mathfrak{D}_{VQ}(d, R_\Delta)$ according to Equation (12) in Equation (16) yields
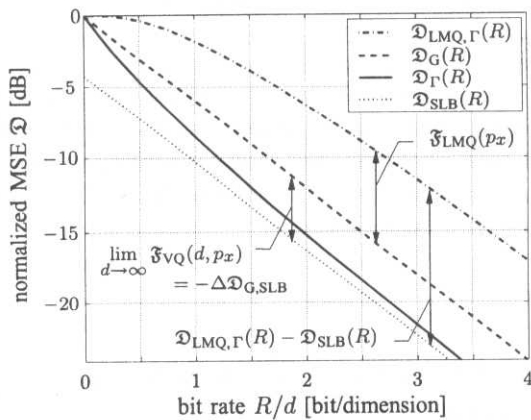


Figure 4. LMQ performance $\mathfrak{D}_{LMQ,\Gamma}(R)$ and distortion-rate function $\mathfrak{D}_\Gamma(R)$ (MSE distortion) for the same Gamma distributed source as in Figure 1.
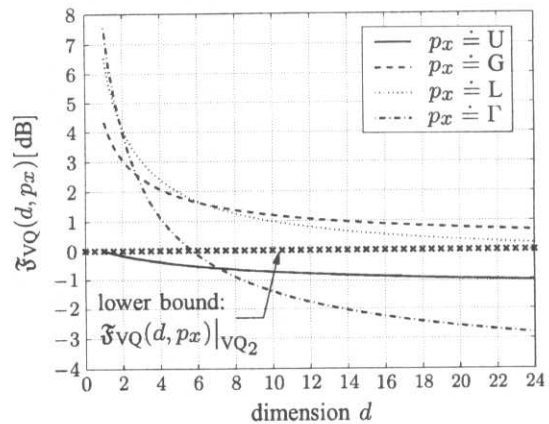


Figure 5. $\mathfrak{F}_{VQ}(d, p_x)$: high-rate MSE distortion of VQ with finite dimension $d \geqslant 1$ relative to the Gaussian DRF for the uniform, Gaussian, Laplacian and Gamma PDF (from Reference [15]).

$$-\frac{6.02}{d}R_\Delta \leqslant -\frac{6.02}{d}R_\Delta + \mathfrak{F}_{VQ}(d, p_{e_2}) \quad [dB] \qquad (19)$$

so that

$$0 \leqslant \mathfrak{F}_{VQ}(d, p_{e_2}) \quad [dB] \qquad (20)$$

Hence, for quantization of $\mathbf{e}_2$ with $VQ_2$, $\mathfrak{F}_{VQ}(d, p_{e_2})$ is bounded below by a value of 0 dB. This implies that, for a given dimension $d$, the shape of $p_{e_2}(e_2)$ must be such as to keep Equation (20) satisfied.

As shown in Figure 5, the $\mathfrak{F}_{VQ}(d, p_{e_2})$-values are required to lie above the 0-dB bound. Note that in case of $p_x$ being a uniform PDF, $p_{e_2}$ can never be uniform since the corresponding $\mathfrak{F}_{VQ}(d, p_{e_2})$-curve lies consistently below 0 dB for a all $d > 1$. For the Gamma PDF, the $\mathfrak{F}_{VQ}(d, p_x)$-curve trespasses the 0-dB bound shortly before $d = 6$, so that for dimensions $d \geqslant 6$ we can be sure that $p_{e_2}$ is neither Gamma nor uniform as both corresponding $\mathfrak{F}_{VQ}(d, p_{e_2})$-curves lie below the 0-dB bound.

(2) *Correlated Sources*: In case of linear dependency or *correlation* $\rho$ between the components of $\mathbf{x}$, we have to consider that correlated sources can be quantized with lower distortion than memoryless sources. Therefore, we can describe $\mathfrak{F}_{VQ}(d, p_x)$ from Equation (13) by the modified quantity $\mathfrak{F}_{VQ}(d, p_x, \rho)$ that depends on the marginal PDF $p_x(x)$ and the correlation $\rho$ according to

$$\mathfrak{F}_{VQ}(d, p_x, \rho) = \mathfrak{F}_{VQ}(d, p_x, \rho = 0) \\ - 10 \log \mathcal{M}(d, \rho) \quad [dB] \qquad (21)$$

where $\mathcal{M}(d, \rho)$ denotes the so-called *memory advantage of VQ* [16]. Under high-rate assumptions $\mathcal{M}(d, \rho)$ coincides with the well-known *spectral flatness measure* of the random variable $x$ [17, 18].

We note that the modified $\mathfrak{F}_{VQ}(d, p_x, \rho)$-curves, compared to the memoryless case of the $\mathfrak{F}_{VQ}(d, p_x)$-curves,[†] are lowered by the respective memory advantage $\mathcal{M}(d, \rho)$. Yet, for the second stage quantizer $VQ_2$ the above $\mathfrak{F}_{VQ}(d, p_x, \rho)$ function is still bounded below as in Equation (20) so that the Gaussian DRF still represents the minimum achievable MSE distortion of $VQ_2$.

For example, with a jointly Gaussian input vector with correlation $\rho$ according to a first-order Markov process, the spectral flatness measure (and hence the memory advantage) can be specified by (e.g. Reference [16])

$$10 \log \mathcal{M}(d, \rho) = d^{-1}(1 - d)10 \log(1 - \rho^2) \quad [dB] \quad (22)$$

[†]Note that in case of a memoryless source with $\rho = 0$, we will, for convenience, retain usage of the notation $\mathfrak{F}_{VQ}(d, p_x)$ rather than $\mathfrak{F}_{VQ}(d, p_x, \rho = 0)$.

According to our earlier observations, the PDF $p_{e_2}$ of the quantization error vector $\mathbf{e}_2$ is allowed to be Gaussian, too. Also, there might still be correlation $\rho_{e_2}$ between the components of $\mathbf{e}_2$. In case of $\rho = 0.5$, and under the assumption that the correlation remains unchanged in the quantization error vector, i.e. $\rho_{e_2} = \rho$, substitution of Equation (22) into (21) shows that for values of $d > 11$, Equation (20) is no longer satisfied, i.e.

$$\mathfrak{F}_{VQ}(d, p_{e_2}, \rho_{e_2})\big|_{p_{e_2} \doteq G, \rho_{e_2} = 0.5} < 0 \, [dB], d > 11 \qquad (23)$$

At first glance, this seems to contradict our previous results which are entirely based on rate-distortion theory. Hence, the only reasonable conclusion that can be drawn here without violating rate-distortion theory is that the amount of correlation is greatly reduced after the first stage of quantization. This means, when the input vector $\mathbf{x}$ is jointly Gaussian with correlation $\rho$ according to a first-order Markov process, the quantization error vector $\mathbf{e}_2$ can also be due to a first-order Gauss-Markov process but with reduced correlation $\rho_{e_2} \ll \rho$. Consequently, the spectral flatness of $\mathbf{e}_2$ and hence the memory advantage $\mathcal{M}(d, \rho_{e_2})$ of $VQ_2$ will be reduced such that

$$0 \leqslant \mathfrak{F}_{VQ}(d, p_{e_2}, \rho_{e_2})\big|_{\rho_{e_2} \ll \rho} \quad [dB] \qquad (24)$$

will be satisfied for all values of $d$.

### 3.2. Upper bound on MSVQ penalty

The development so far provided us with some fundamental insight into the performance of MSVQ. While the first quantizer stage of MSVQ performs like conventional fixed rate VQ, the performance of the following quantizer stages depends on the properties of the residual after the first stage.

For MSVQ with $i = 1, \ldots, \nu$ stages, let $R_i$ denote the bit rate given in bit/vector that is allocated to the individual quantizer stages $VQ_i$. It is convenient to specify the maximum MSE distortion after $\nu$ stages relative to the minimum attainable MSE distortion $\mathfrak{D}_{VQ}(d, R_\Sigma)$ of a singlestage $VQ_\Sigma$ with rate

$$R_\Sigma = \sum_{i=1}^{\nu} R_i \qquad (25)$$

This way, we will obtain an upper bound on the *performance penalty* of MSVQ.

Considering only the first stage, the performance penalty is trivially 0 dB as $R_\Sigma = R_1$. Given that $R_1$ is sufficiently large so that $VQ_\Sigma$ has already reached asymptotic

behaviour, a further increase of $R_\Sigma$ leads to a further decrease of $\mathfrak{D}_{\mathrm{VQ}}(d, R_\Sigma)$ at the asymptotic rate of 6.02/$d$ dB/bit. Consequently, the performance penalty of each individual quantizer stage beyond the first one, i.e. for $i \geqslant 2$, is given in decibels by the corresponding quantity $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i}, \rho_{e_i})$ from Equation (21), where $p_{e_i}$ is the marginal PDF of the vector $\mathbf{e}_i$ which is the input to the $i$th quantizer stage $\mathrm{VQ}_i$ respectively the quantization error of $\mathrm{VQ}_{i-1}$. Therefore, the overall performance penalty $\Delta\mathfrak{D}_{\mathrm{MSVQ}}(d, p_x, \nu)$ after $\nu$ stages is obtained by

$$\Delta\mathfrak{D}_{\mathrm{MSVQ}}(d, p_x, \nu) = \sum_{i=2}^{\nu} \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i}, \rho_{e_i}) \quad [\mathrm{dB}] \quad (26)$$

To specify the maximum $\Delta\mathfrak{D}_{\mathrm{MSVQ}}(d, p_x, \nu)$, we have to make some worst case assumptions on the shape of $p_{e_i}$ and the correlation $\rho_{e_i}$ between the components of $\mathbf{e}_i$, given the marginal PDF $p_x$ and the correlation $\rho$ of the input vector $\mathbf{x}$, so that the $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i}, \rho_{e_i})$ terms in the above Equation (26) take on their respective maximum values.

For the jointly Gaussian case with correlation $\rho$ according to a first-order Markov process, only the maximum correlation $\rho_{e_i}^{\max}(d)$ that potentially exists in $\mathbf{e}_i$ can be calculated by solving Equation (24) for $\rho_{e_i}$. With Equations (22) and (21) we obtain

$$\rho_{e_i}(d) \leqslant \rho_{e_i}^{\max}(d) = \left[1 - \left(\frac{\epsilon^2(p_{e_i})\big|_{p_{e_i} \doteq \mathrm{G}}}{\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_{i-1}})}\right)^{\frac{1}{1-d}}\right]^{\frac{1}{2}} \quad (27)$$

The only worst case assumption that is admissible here is thus, to consider all quantization error vectors $\mathbf{e}_i$ as memoryless, which indeed represents the ultimate worst case assumption with respect to the memory advantage of the following quantizer stages. With the above assumption, the $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i}, \rho_{e_i})$ terms in Equation (26) reduce to their memoryless form $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i}, \rho_{e_i} = 0) \equiv \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$. The potential maximum values of $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$ are due to a worst case assumption on the respective shape of $p_{e_i}$. As will be shown by a reinspection of Figure 5, the worst case PDF $p_{e_i}$ individually depends on the shape of $p_{e_{i-1}}$ being the marginal PDF of the input vector to the $(i-1)$th quantizer stage. Hence, maximum values of $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$ imply maximization with respect to $p_{e_{i-1}}(e_{i-1})$, thus for $i = 2, \ldots, \nu$

$$\max_{p_{e_{i-1}}}\{\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\} = \mathfrak{F}_{\mathrm{VQ}}^{\max}(d, p_{e_i} | p_{e_{i-1}}) \quad [\mathrm{dB}] \quad (28)$$

whereas $p_{e_1} \equiv p_x$.

We conjecture that under high-rate assumptions the PDF of the quantization error resulting from PDF-optimized

VQ tends to be less sharp and less peaked than that of the input to the quantizer. This statement is also supported by our observation from the $\mathfrak{F}_{\mathrm{VQ}}(d, p_x)$-curves in Figure 5, which indicate that for increasing dimension $d$ the shape of the marginal PDF of the quantization error becomes increasingly similar to the smooth shape of the Gaussian PDF. Analog, as $d \to \infty$ the values of $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$ converge to those values for the Gaussian PDF.

Figure 6 shows a magnified display of the $\mathfrak{F}_{\mathrm{VQ}}(d, p_x)$-curves for small dimensions in the range of $2 \leqslant d \leqslant 4$. In case of $p_{e_{i-1}}$ being, for example, Gamma the intersection of the respective $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$-curve with the $\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})$-curve for the Gaussian PDF lies between $2 < d < 3$. As for $d = 2$,

$$\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_{e_i} \doteq \Gamma} > \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_x \doteq \mathrm{G}} \quad [\mathrm{dB}]$$

the Gamma PDF represents the worst case PDF $p_{e_i}$, so that

$$\mathfrak{F}_{\mathrm{VQ}}^{\max}(d, p_{e_i} | p_{e_{i-1}}) = \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_{e_i} \doteq \Gamma} \quad [\mathrm{dB}]$$

For $d = 4$ we have,

$$\mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_{e_i} \doteq \Gamma} < \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_x \doteq \mathrm{G}} \quad [\mathrm{dB}]$$

so that the Gaussian PDF represents the worst case PDF $p_{e_i}$, and

$$\mathfrak{F}_{\mathrm{VQ}}^{\max}(d, p_{e_i} | p_{e_{i-1}}) = \mathfrak{F}_{\mathrm{VQ}}(d, p_{e_i})\big|_{p_{e_i} \doteq \mathrm{G}} \quad [\mathrm{dB}]$$

Note that in case of $p_{e_{i-1}}$ being uniform, the Gaussian PDF represents the worst case PDF $p_{e_i}$ for all dimensions $d > 1$, as the $\mathfrak{F}_{\mathrm{VQ}}(d, p_x)$-curve of the uniform PDF lies consistently below 0 dB and thus below the $\mathfrak{F}_{\mathrm{VQ}}(d, p_x)$-curve for $p_{e_i}$ being Gaussian.
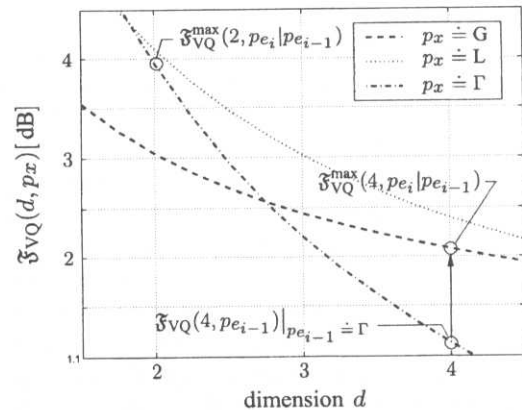


Figure 6. Magnified display of the $\mathfrak{F}_{\mathrm{VQ}}(d, p_x)$-curves from Figure 5 in the range of $2 < d < 4$.

Generally, the maximum value of $\mathfrak{F}_{\mathrm{VQ}}(d,p_{e_i})$ for arbitrary shapes of $p_{e_{i-1}}$ is given by

$$\mathfrak{F}_{\mathrm{VQ}}^{\max}(d,p_{e_i}\,|\,p_{e_{i-1}}) = \mathfrak{F}_{\mathrm{VQ}}(d,p_{e_{i-1}}) \quad [\mathrm{dB}] \qquad (29)$$

under the constraint of

$$\mathfrak{F}_{\mathrm{VQ}}(d,p_{e_{i-1}}) \geqslant \mathfrak{F}_{\mathrm{VQ}}(d,p_{e_i})\big|_{p_x \doteq \mathrm{G}} \quad [\mathrm{dB}] \qquad (30)$$

Otherwise, $\mathfrak{F}_{\mathrm{VQ}}^{\max}(d,p_{e_i}\,|\,p_{e_{i-1}})$ is given by the respective value for the Gaussian PDF, i.e.

$$\mathfrak{F}_{\mathrm{VQ}}^{\max}(d,p_{e_i}\,|\,p_{e_{i-1}}) = \mathfrak{F}_{\mathrm{VQ}}(d,p_{e_i})\big|_{p_{e_i} \doteq \mathrm{G}} \quad [\mathrm{dB}] \qquad (31)$$

These $\mathfrak{F}_{\mathrm{VQ}}^{\max}(d,p_{e_i}\,|\,p_{e_{i-1}})$ functions specify individual upper bounds on the performance penalty for each quantizer stage $i$. With the expression for $\Delta\mathfrak{D}_{\mathrm{MSVQ}}(d,p_x,\nu)$ in Equation (26), the desired upper bound on the overall performance penalty of MSVQ with $\nu$ quantizer stages is given by

$$\Delta\mathfrak{D}_{\mathrm{MSVQ}}^{\max}(d,p_x,\nu) = \sum_{i=2}^{\nu} \mathfrak{F}_{\mathrm{VQ}}^{\max}(d,p_{e_i}\,|\,p_{e_{i-1}}) \quad [\mathrm{dB}] \qquad (32)$$
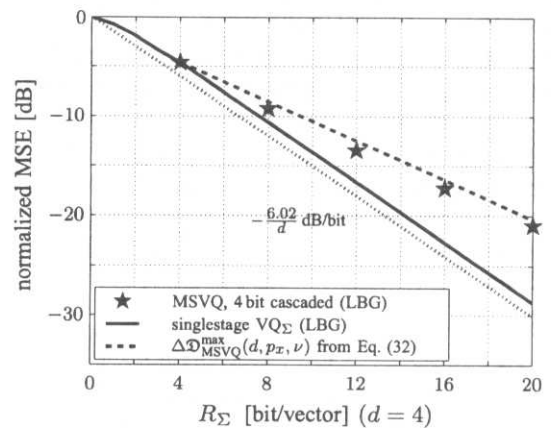
## 4. EXPERIMENTAL RESULTS

The specification of $\Delta\mathfrak{D}_{\mathrm{MSVQ}}^{\max}(d,p_x,\nu)$ has been entirely devoted to results on theoretical VQ performance based on high rate assumptions. To see how this bound applies in practice and for relatively small rates, MSVQ is applied to a first-order Gauss–Markov source.

Such Gauss–Markov source with adjacent correlation $\rho$ can be used for simplified AR(1) modeling of speech (comp. Reference [10]). Therefore, an i.i.d. zero-mean Gaussian random variable $z$ with variance $\sigma_z^2 = 1$ was used to excite a first-order, linear recursive filter. The filter output $x$ was partitioned into a set of 700 000 training vectors $x$ with dimension $d = 4$. The vector quantizers were individually optimized to the PDF of the respective training set by the LBG algorithm [11].
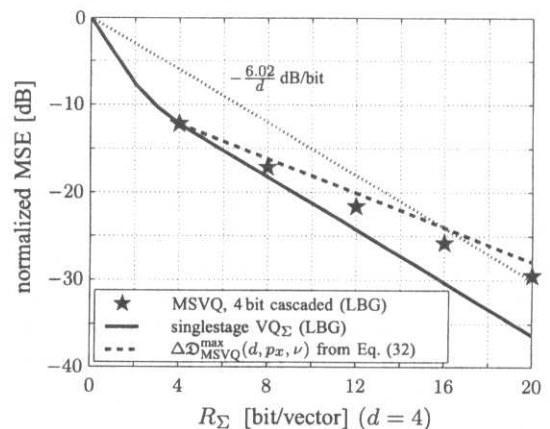
The experiment was carried out for (a) a memoryless source with $\rho = 0$ and (b) a highly correlated source with $\rho = 0.95$. For these two settings, the plots in Figure 7 (a) and (b) show both the predicted as well as the actually measured MSE distortion as a function of the bit rate $R_\Sigma$ respectively. In both figures, the solid line shows the performance of singlestage $\mathrm{VQ}_\Sigma$ which constitutes the minimum attainable MSE distortion $\mathfrak{D}_{\mathrm{VQ}}(d,R_\Sigma)$ for quantization of that particular source, while the stars indicate the performance of MSVQ with

$\nu = 5$ quantizer stages $\mathrm{VQ}_i$, each employing $R_i = const. = 4\,\mathrm{bit/vector}$.

The plots in Figure 7(a) and (b) clearly indicate, that at least for a cascade of more than two VQs the performance of MSVQ is significantly worse than that of singlestage $\mathrm{VQ}_\Sigma$. For quantization of the memoryless Gaussian source it is now interesting to compare the respective performances of MSVQ and LMQ. Thus for a given bit rate of $R_\Sigma$ bit/vector, the MSE distortion of MSVQ from Figure 7 (a) is compared to the MSE distortion of LMQ as shown by the dashed line in Figure 4 given the corresponding bit rate of $R = R_\Sigma/d$ bit/sample. The surprising result is that already from the third quantizer stage $\mathrm{VQ}_3$ the MSE



(a) Gauss-Markov source: $\rho = 0.0$



(b) Gauss-Markov source: $\rho = 0.95$

Figure 7. Simulation results for MSE distortion obtained after each stage of MSVQ versus MSE distortion of singlestage $\mathrm{VQ}_\Sigma$. For comparison $\Delta\mathfrak{D}_{\mathrm{MSVQ}}^{\max}(d,p_x,\nu)$ specifies a new upper bound on the MSE distortion after $\nu$ stages of MSVQ according to Equation (32).

distortion of MSVQ is higher than that of LMQ, which implies that for a memoryless Gaussian source a cascade of $\nu \geqslant 3$ VQ stages performs actually worse than LMQ.

Regarding the performance of the singlestage $VQ_\Sigma$, both $\mathfrak{D}_{VQ}(d, R_\Sigma)$-curves from Figure 7(a) and (b) reach their asymptotic behaviour of $-6.02/d$ dB/bit at $R_\Sigma = 5$ bit/vector. A comparison of the $\mathfrak{D}_{VQ}(d, R_\Sigma)$-curves further indicates that VQ of the correlated source (comp. Figure 7(b)) yields lower MSE distortion than VQ of the memoryless source (comp. Figure 7(a)) for all bit rates $R_\Sigma$ which complies with our theoretical considerations on the memory advantage of VQ according to Equation (21). Between these two curves we measure an asymptotic difference of $10 \log \widetilde{\mathcal{M}}(d, \rho) = 7.58$ dB. Note that this measured value of the memory advantage is exactly the same as what can be theoretically calculated from Equation (22), i.e. $\widetilde{\mathcal{M}}(d, \rho) = \mathcal{M}(d, \rho)$.

Due to the worst case assumption of zero correlation $\rho_{e_i}$ in the quantization error vectors $\mathbf{e}_i$, the predicted MSVQ penalty $\Delta \mathfrak{D}_{MSVQ}^{max}(d, p_x, \nu)$ must be the same for both cases, which can be readily recognized from the identical slope of the dashed lines in Figure 7(a) and (b). It is therefore, easy to understand that the bound is more accurate for the memoryless case than for $\rho = 0.95$. In the latter case the actual correlation values $\rho_{e_i}$ are non-zero which causes a considerable memory advantage for each $VQ_i$ beyond $VQ_1$. Column 2 of Table 3 shows the individual memory advantage for each $VQ_i$ as can be determined from the distance between the actually measured values of MSVQ distortion (stars) in Figure 7(b) and (a) minus the offset of 7.58 dB (see above). Insertion of these values into Equation (22) and some algebra yields an expression to specify the effective $\rho_{e_i}$ values as listed in column 3.

First, we note the rapid decrease of correlation after the first stage as indicated by the differential value of $\Delta \rho_2 = 0.65$ in column 4, while the further decrease of $\rho_{e_i}$ is comparably flat. Secondly, we note a value of

$\rho_{e_2} = 0.3$. This value is considerably lower than the potential maximum value of $\rho_{e_2}^{max} = 0.69$ which can be calculated from Equation (27). Together with the results from the above performance comparison with LMQ the results on the effective $\rho_{e_i}$ values after the first quantizer stage may explain, why in practice MSVQ schemes usually have no more than two or three stages.

## 5. CONCLUSIONS

In this paper, we addressed the performance penalty that arises when using a hierarchical coder based on MSVQ as opposed to the use of a non-hierarchical coder based on fixed rate VQ when using the MSE distortion criterion. Therefore, we reviewed several important results from rate-distortion and high-rate quantization theory. The theoretical tools that these two theories provide for describing the performance of fixed-rate VQ were used to specify a new upper bound on the performance penalty of MSVQ. From the theoretical analysis of MSVQ we also gained fundamental insight in the performance of hierarchical coding based on cascaded coding approaches. The results on MSVQ performance respectively the MSVQ penalty were confirmed by practical experiments. These experiments also indicated that the upper bound on the MSVQ penalty, although completely based on theoretical considerations and high-rate assumptions, applies for low bit rates as well.

Table 3. Experimental results for MSVQ of a Gauss–Markov source with adjacent correlation $\rho = 0.95$.

| $VQ_i$ | $\widetilde{\mathcal{M}}(d, \rho)\big|_{VQ_i}$ [dB] | $\rho_{e_i}$ | $\Delta \rho_i$ |
|---|---|---|---|
| 1 | 7.58 | 0.95 | |
| 2 | 0.31 | 0.30 | 0.65 |
| 3 | 0.25 | 0.28 | 0.02 |
| 4 | 0.33 | 0.31 | −0.03 |
| 5 | 0.03 | 0.10 | 0.21 |

## REFERENCES

1. Berger T. *Rate Distortion Theory*. Prentice Hall: Englewood Cliffs, New Jersey, 1971.
2. Gray RM. *Source Coding Theory*. Kluwer Academic Publishers: Boston, USA, 1990.
3. Gersho A, Gray R. *Vector Quantization and Signal Compression*. Kluwer Academic Publishers: Boston, USA, 1992.
4. Cover TM, Thomas JA. *Elements of Information Theory*. Wiley: New York, 1991.
5. Noll P, Zelinski R. Bounds on quantizer performance in the low bit rate region. *IEEE Transactions on Communications* 1978; **COM-26**(2):300–304.
6. Blahut R. Computation of channel capacity and rate distortion functions. *IEEE Transactions on Information Theory* 1972; **IT-18**: 460–473.
7. Lloyd SP. Least squares quantization in PCM. *IEEE Transactions on Information Theory* 1982; **28**(2):129–137.
8. Max J. Quantizing for minimum distortion. *IRE Transactions on Information Theory* 1960; **6**:7–12.
9. Algazi V. Useful approximations to optimum quantization. *IEEE Transactions on Communications* 1966; **COM-14**:297–301.
10. Jayant NS, Noll P. *Digital Coding of Waveforms*. Prentice Hall: Englewood Cliffs, New Jersey, USA, 1984.

11. Linde Y, Buzo A, Gray R. An algorithm for vector quantizer design. *IEEE Transactions on Communications* 1980; **28**(1):84–95.

12. Zador P. Asymptotic quantization error of continuous signals and the quantization dimension. *IEEE Transactions on Information Theory* 1982; **IT-28**(2):139–149.

13. Gersho A. Asymptotically optimal block quantization. *IEEE Transactions on Information Theory* 1979; **IT-25**(4):373–380.

14. Conway J, Sloane N. Voronoi regions of lattices, second moments of polytopes, and quantization. *IEEE Transactions on Information Theory* 1982; **IT-28**(2):211–226.

15. Erdmann C. Hierarchical vector quantization: theory and application to speech coding. Ph.D. dissertation, Aachener Beiträge zu digitalen Nachrichtensystemen, ed. P. Vary, Bd. 18, RWTH Aachen, 2002.

16. Lookabaugh T, Gray R. High resolution quantization theory and the vector quantizer advantage. *IEEE Transactions on Information Theory* 1989; **35**(5):1020–1033.

17. Markel JD, Gray AH, Jr. *Linear Prediction of Speech*. Springer, 1976.

18. Makhoul J, Roucos S, Gish H. Vector quantization in speech coding. *Proceedings of the IEEE* 1985; **73**(11):1551–1588.

## AUTHORS' BIOGRAPHIES

**Christoph Erdmann** received the Dipl.-Ing. degree in electrical engineering in 1998 from Aachen University of Technology, Germany. He has, since then, been with the Institute of Communication Systems and Data Processing, working on speech coding algorithms for heterogeneous communication networks with a focus on hierarchical coding techniques. He is currently pursuing his Ph.D. with the thesis, *Hierarchical Vector Quantization: Theory and Application to Speech Coding*.

**Peter Vary** received the Dipl.-Ing. degree in electrical engineering in 1972 from the University of Darmstadt, Germany. In 1978, he received his Ph.D. from the University of Erlangen-Nuremberg and in 1980 he joined Philips Communication Industries (PKI) in Nuremberg, Germany. He became head of the Digital Signal Processing Group, which made substantial contributions to the development of GSM. Since 1988 he has been a Professor at Aachen University of Technology, Germany, and head of the Institute of Communication Systems and Data Processing. His main research interests are speech coding, channel coding, error concealment, adaptive filtering for acoustic echo cancellation and noise reduction and concepts of mobile radio transmission.