

A Trellis-Based Optimal Parameter Value Selection for Audio Coding

Ashish Aggarwal, *Member, IEEE*, Shankar L. Regunathan, *Member, IEEE*, and Kenneth Rose, *Fellow, IEEE*

Abstract—This paper considers the problem of selecting a set of parameter values from a given parameter space, in order to perform rate-distortion optimization in the context of audio compression. Due to interdependencies between parameters, separate optimization of parameter values is inherently suboptimal, yet a straightforward brute-force joint search involves prohibitive computational complexity. This work proposes a new method for joint rate-distortion optimization, while accounting for interparameter dependencies. The optimal solution is achieved, at significantly reduced complexity as compared to a brute-force search, by employing a Viterbi search over a trellis. Two objective distortion metrics are specifically considered: the average, and the maximum noise-to-mask ratio. Subjective (AB/MOS) and objective (average/maximum noise-to-mask ratio) tests demonstrate considerable gains at low bit rates of 16 kbps per channel for a 44.1-kHz sampled audio signal using the proposed approach.

Index Terms—Advanced audio coder (AAC), audio coding, bit allocation, dynamic programming, parameter selection, side-information, trellis, Viterbi.

I. INTRODUCTION

AUDIO COMPRESSION is central to many multimedia applications such as digital audio broadcasting and transmission of music over the Internet. Such applications benefit substantially from improved compression performance. Current audio coders such as MPEG Advanced Audio Coder (AAC) [1], [2], AC3 [3], PAC [4], ATRAC [5], and G.722.1 [6] rely heavily on the removal of perceptually irrelevant information [7]–[10] from the source signal. For a thorough description of current audio coding techniques, see [11]. Perceptually irrelevant information is exploited via calculation of the masking threshold—the threshold below which a signal (or noise) is rendered inaudible—which, in turn, involves time-adaptive spectral shaping of the quantization noise. Shaping of the quantization noise is a rate-distortion optimization performed at the encoder. Noise shaping is typically achieved by varying the granularity of the quantizer employed in the different

frequency bands (or critical bands [7] that emulate the human auditory system’s grouping of adjacent frequency bands). The choice of quantizer granularity is one of the many parameters whose values are chosen dynamically by the encoder in order to perform rate-distortion optimization. We refer to the complete set of such parameters as the “encoding parameters.” Selection of encoding parameter values is central to the rate-distortion optimization performed by the encoder.

Consider AAC for example. It performs spectral decomposition of a frame of the audio signal, groups the spectral coefficients into bands, and quantizes the coefficients using scalar quantizers. Adaptive noise shaping is achieved by allowing per-band scaling of the generic scalar quantizer by an appropriate scale factor (SF). Since the SF is shared by the entire band, each band is commonly referred to as a scale factor band (SFB). The quantized coefficient indices are entropy coded using a possibly different Huffman codebook (HCB) for each SFB. The choice of the HCB is made from a set of predesigned codebooks. The SF and HCB values chosen per SFB form the set of parameters which, together with the quantized coefficient indices, convey to the decoder all the information needed to reconstruct the coefficients for the frame. These parameters constitute the *encoding parameters*, whose values are determined by the encoder for every frame of the audio signal. It is conceivable to obtain the optimal parameter values in a rate-distortion sense using a straightforward brute-force search. However, such an optimal scheme involves prohibitive computational complexity due to the large size of the parameter space. AAC allows for as many as 60 distinct SF values and 12 predesigned HCBs. For a frame of 44.1-kHz sampled audio consisting of 49 SFBs the cardinality of the parameter space reaches $(60 \times 12)^{49}$ —clearly putting brute-force search beyond computational reach.

A suboptimal choice of parameter values can significantly degrade the encoder’s compression performance. At relatively high encoding rates, there exist multiple solutions for which the quantization noise completely falls below the masking threshold. In this case, a suboptimal choice in the rate-distortion sense may not cause considerable subjective performance degradation. However, when the signal is quantized at low rates (for example, 16–48 kbps/channel for a 44.1-kHz sampled signal) it is impossible to maintain all the quantization noise below the masking threshold, and it is critical to carefully optimize the parameter values. Hence, computationally efficient search for the optimal encoding parameter values is an interesting and important problem in audio coding. It is known to play a crucial role in other signal compression applications as well [12]–[14]. In this paper we focus on the

Manuscript received January 12, 2003; revised November 24, 2004. This work was supported in part by the NSF under Grants MIP-9707764, EIA-9986057, and EIA-0080134, the University of California MICRO Program, Dolby Laboratories, Inc., Lucent Technologies, Inc., Mindspeed Technologies, and Qualcomm, Inc. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Ravi P. Ramchandran.

A. Aggarwal is with Harman Consumer Group, Northridge, CA 93129 USA (e-mail: aaggarwa@harman.com).

S. L. Regunathan is with Microsoft Corp., Redmond, WA 98052 USA (e-mail: shrane@microsoft.com).

K. Rose is with the Department of Electrical and Computer Engineering, University of California, Santa Barbara, CA 93106-9560 USA (e-mail: rose@ece.ucsb.edu).

Digital Object Identifier 10.1109/TSA.2005.855833

problem of optimally selecting encoding parameter values for audio compression. The term “optimality” is employed in the rate-distortion sense, i.e., the optimal selection is one which minimizes the distortion measure for the prescribed total rate. We outline the solution for two objective metrics: the average and the maximum noise-to-mask ratios (NMR) [15]–[18]. Note that this paper does not directly address the widely recognized problem of finding an objective metric that adequately reflects the subjective quality of reconstructed audio signals.

Selection of values for the encoding parameters is closely related to the problem of bit allocation [14], [19] whose early approaches employ high-resolution quantization theory to arrive at a simple solution that is implementable by the popular water-filling algorithm [20]–[22]. The algorithm attempts to maintain a constant distortion (say, NMR) across the coefficients (or critical bands) and forms the basis for selection of parameter values in various audio coding algorithms, such as the so-called two-loop search (TLS) [23]. For a comprehensive review of approaches to bit allocation and parameter value selection, see [12], [13], [24]–[27]. Conventional water-filling based approaches suffer from two major drawbacks. First, the coefficients are not statistically independent, however, conventional methods do not accurately account for these inter-coefficient dependencies that exist in the spectral representation [13]; and second, as we show later, their solution fails to distinguish between the objective measures considered in this paper. Consequently, the choice of encoding parameter values may be significantly suboptimal for either metric and the resulting compression performance penalty may be considerable at low bit rates. The importance of improved low bit rate performance is further highlighted in the case of bit rate scalable (also, embedded or layered) compression [28], [29] where multiple low bit rate encoding modules are employed.

We propose a search algorithm which explicitly optimizes for the interparameter dependencies that exist in the spectral representation. To combat the prohibitive computational complexity of the straightforward brute-force solution, we recast the problem as a search through a trellis, and employ dynamic programming [30] to obtain the optimal solution at a drastically reduced search complexity. The search is outlined for the two objective metrics, which are both based on the NMR, namely, ANMR and MNMR. The proposed trellis-based search is compared with the water-filling approach of TLS described in [23] as competing search modules in AAC (see Section V for further details). Note that TLS is the best publicly disclosed search method for AAC. Simulation results demonstrate substantial improvement in the encoder’s low bit rate performance. For example, on a standard critical test database from EBU-SQAM [31], [32] comprising of 44.1-kHz sampled (mono) audio signal, the proposed search method operating at bit rates in the range of 16–32 kbps, requires half the bit rate to achieve the same objective (ANMR/MNMR) and subjective (AB/MOS) quality as TLS. When implemented within a four-layer scalable coder where each layer employs 16-kbps AAC encoding modules, the proposed scheme achieved performance close to that of a 56 kbps non-scalable AAC coder. Furthermore, as the solution achieves rate-distortion optimality, it promises a useful framework for performance evaluation of other search schemes (e.g.,

see [33] and [34]). The performance benefit is achieved at the expense of computational complexity as compared to the TLS and it is incurred only at the encoder. It is important to emphasize that the proposed scheme leaves the bit stream syntax intact and the AAC decoder unaltered. The method is hence standard-compatible. Preliminary results of this work have been reported in [35] and [36].

The organization of the paper is as follows. Section II provides a brief background to the problem. The proposed trellis-based search method is derived in Section III. The implementation of the proposed search within AAC is described in Section IV, and results are summarized in Section V.

II. BACKGROUND

A. Objective Measures in Audio Coding

Most objective measures employed in rate-distortion optimization of the encoder are designed to model subjective, perceptual distortion. On the one hand, simple metrics such as the mean-squared error (MSE) fail to model perceptual distortion accurately. On the other, metrics with relatively good modeling accuracy, such as PAQM [37] and PEAQ [38], [39], are too complex to be used in run-time optimization of the encoder. While a suitable objective metric that accurately models the subjective quality remains an unsolved problem, most widely used objective measures involve the NMR [15], [16], which is the ratio of the quantization noise energy to the masking threshold in the given critical band [7]–[10]. The NMR in the critical band may equivalently be viewed as a weighted squared error (WSE) whereby the weights are simply the inverse of the masking threshold in the critical band. NMR below unity in a critical band indicates that quantization noise in that band is imperceptible. At low rates it is often impossible to maintain the NMR below unity in all the critical bands. Hence, the NMR values obtained from the various critical bands are combined into a scalar distortion metric. Two common metrics are: ANMR, which is the NMR averaged over all the critical bands in the frame, and MNMR, which is the maximum NMR of all the critical bands in a frame [17], [18].

Let d_i be the squared quantization error, w_i be the weight of critical band i , and N be the total number of bands. ANMR is given by

$$D_A = \frac{1}{N} \sum_{i=1}^N w_i d_i \quad (1)$$

and MNMR by

$$D_M = \max_{i=1}^N w_i d_i. \quad (2)$$

Subjective listening tests performed by us [36] and in [17] and [18] indicate substantial differences in the quality of audio signal resulting from optimization of the two metrics. At low rates, optimization of MNMR metric resulted in fewer annoying artifacts such as clicks, but the average quality was perceived to be inferior to ANMR. However, there was no general consistent preference for either.

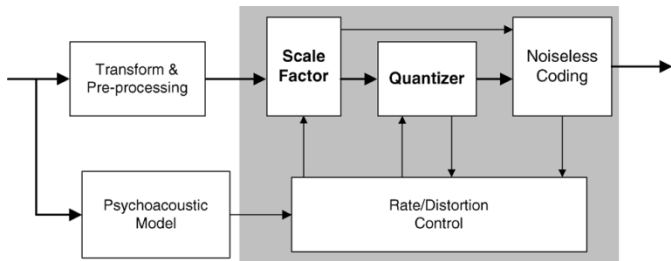


Fig. 1. Block diagram of the AAC encoder. Transform and preprocessing tools are applied prior to quantization and coding (QC). The psychoacoustic model outputs the masking threshold which is used for rate-distortion optimization.

B. MPEGs Advanced Audio Coding

This section focuses on the quantization module of AAC. A simplified, high-level block diagram of the AAC encoder is shown in Fig. 1. The quantization and coding (QC) module, which is central to this work, is shown in greater detail. The time domain signal is grouped into overlapping frames and transformed into the spectral domain using the modified discrete cosine transform (MDCT). The transform yields a set of 1024 coefficients that are then quantized using the QC module. In the QC module, the transform coefficients are grouped into nonuniform frequency bands, termed scale factor band (SFB), and all coefficients within a given SFB are quantized using the same nonuniform scalar quantizer which is characterized using a compander (see [1] and [2], for further details). The quantizer is a scaled version of the generic quantizer, and is determined by the scale factor (SF) parameter, which is selected for each SFB and controls the desired noise level in the band. The time domain signal is also input to the psychoacoustic model, whose output is the masking threshold for each SFB.

Statistical redundancy in the quantized coefficient indices is exploited by the use of entropy and run-length coding techniques. AAC offers a set of 12 predesigned Huffman codebooks (HCB), from which one is selected for each SFB for encoding the quantized coefficient indices. In addition to the quantized coefficient indices, side information must be transmitted to specify SF and HCB selections for each SFB. SF values are differentially encoded using a variable length code, and HCB selection is encoded using a run-length code. The rate-distortion optimization at the encoder involves the choice of SF and HCB values for each SFB.

C. Parameter Value Selection in Current Audio Encoders

Recall that removal of perceptually irrelevant information via quantization noise shaping is implemented in audio coding by appropriately selecting the values of the encoding parameters for the various frequency bands. This problem is, in turn, closely related to the problem of bit allocation, which has been extensively covered in the signal compression literature. A comprehensive coverage of this topic is beyond the scope of this paper, and can be found in [14] and [19]. We will only briefly outline here the relevant portions of the classic problem of bit allocation and its known water-filling solution [20]–[22] which stems from high-resolution quantization theory.

The bit allocation problem is one where a fixed bit budget needs to be distributed among different coefficients in order

to minimize the distortion (e.g., NMR) at hand. Let b_i be the number of bits allocated to, and d_i be the resulting distortion of, coefficient i . Let R_t be the target rate and N be the total number of coefficients in the frame. The problem of bit allocation may be stated as

$$\mathbf{b}^* = \arg \min_{\mathbf{b}: \sum_{i=1}^N b_i \leq R_t} \sum_{i=1}^N d_i(b_i) \quad (3)$$

where $\mathbf{b} = (b_1, \dots, b_N)$ represents the bit allocation vector and \mathbf{b}^* is the optimal allocation. Early solutions to the problem of bit allocation use high-resolution (quantization) approximation [20]–[22] to model the distortion as

$$d_i(b) = k\nu_i^2 2^{-2b} \quad (4)$$

where ν_i^2 is the variance of coefficient i , and k is a constant that depends on the slope of the probability density function of the coefficients. All coefficients are typically assumed to have the same k . This model is the basis of the celebrated solution to the problem of independent bit allocation

$$b_i^* = \frac{R_t}{N} + \frac{1}{2} \log_2 \frac{\nu_i^2}{\omega^2} \quad (5)$$

where $\omega^2 = (\prod_{i=1}^N \nu_i^2)^{1/N}$ (for proof see [14]). When the bit allocation is optimal, it is easy to see that the resulting distortion is the same for all coefficients, i.e.,

$$d_i(b_i^*) = k\omega^2 2^{-\frac{2R_t}{N}}. \quad (6)$$

Hence, the optimal bit allocation can be implemented by a simple water-filling algorithm, where the same level of distortion is maintained at all coefficients, and this level is varied to meet the target rate. Note that in the context of audio, (3) corresponds to the ANMR measure. An interesting (and perhaps surprising) observation is made when one analyzes the bit allocation problem for minimizing the MNMR distortion metric. It turns out that the same water-filling solution optimizes MNMR metric as well, *at high resolution* (see the Appendix for details).

Variants of the basic water-filling algorithm are typically employed for selection of parameter values in audio coding. Consider, for example, TLS [23], which consists of two nested loops. The task of the inner iteration loop is to uniformly change the SF values of all the SFBs by a constant amount, and determine the HCB values so that the given spectral data may be encoded while satisfying the rate constraint. The outer loop changes the SF values of individual SFBs, and thus shapes the quantization noise to best match the psychoacoustic model. In a nutshell, TLS tries to maintain the NMR in each SFB below a given level, and then adjusts this level to meet the rate constraint.

One major drawback of the approach is the use of the distortion model given in (4). The model makes it difficult, and often impossible, to account for the side-information rate when performing dynamic bit allocation. Shoham and Gersho proposed an alternative Lagrangian-based solution to account for the side-information rate [13], [40], without recourse to high-resolution approximation or other analytical models of the distortion. However, they assumed coefficient independence in

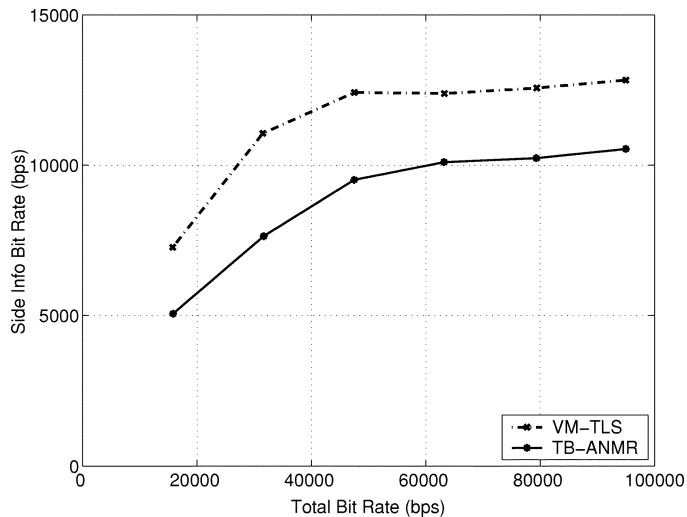


Fig. 2. Side-information employed by the TLS for a AAC implementation using VM-TLS and TB-ANMR (proposed). The side-information rate is plotted versus the total rate for a single channel 44.1-kHz sampled audio signal. Side-information includes bits consumed to transmit SF and HCB values.

calculating the side-information rate. Similar results were also reported in [24] and [25]. The more general case of dependent bit allocation was addressed in [12].

D. Problem Motivation and Challenges

The encoder's problem is to select the values of the encoding parameters so as to minimize the distortion metric for the given target rate. This problem is complicated by several factors. As the statistical characteristics of the audio signal vary considerably with time, parameter values must be chosen dynamically. A trade-off emerges wherein dynamic selection helps reduce the rate required to transmit the quantized coefficients but must be transmitted as side-information and hence increases the rate. Further, there exist dependencies across the spectral coefficients (or critical bands) which affect the total bit rate. These dependencies are, in fact, the motivation behind AACs use of run-length and differential coding of HCB and SF values, respectively. Thus, the side-information rate (and hence the total rate) is a joint function of all parameter values used to encode the coefficients in the frame. It cannot be expressed as a simple sum of the bits independently optimized for encoding individual parameter values. This observation points to a major shortcoming of the conventional water-filling approach, which relies critically on the invalid assumption of parameter independence. Yet another drawback of the conventional approach is due to the underlying rate-distortion model, which is derived from high-resolution quantization theory. The model not only breaks down when encoding rates are low, but also fails to accurately account for the (time varying) rate required to transmit the side-information. Conventional schemes do not take parameter dependencies into account and fail to explicitly optimize the side-information rate. TLS, in particular, accounts for the side-information rate only by counting the side-information bits in the inner (rate) loop. However, it does not explicitly optimize the encoding parameter values while accounting for their contribution to the side-information rate. At high rates, the price of ignoring explicit optimization of the side-information rate may be tolerable because the

side-information rate forms a relatively small percentage of the total rate. Fig. 2 shows the rate consumed in transmission of SF and HCB values versus the total rate. It is evident that, at low bit rates, side-information may consume as much as 30%–40% of the total rate. At these rates, ignoring side-information and parameter dependencies often results in a severe performance penalty.

The problem is further complicated in the case of audio by the fact that different objective criteria, such as ANMR and MNMR, may be used for encoder optimization. Note that this complication disappears whenever the assumptions of high-resolution and parameter independence are valid. Recall further that in this case the same water-filling algorithm optimizes both criteria. Thus, the TLS-based search method can afford to be agnostic of the distortion metric. However, these assumptions fail to hold in practical audio coding. In fact, subjective tests [17], [18], [36] indicate that in practice (when the above assumptions do not hold) the perceived output quality of the optimum solution for the two measures differ significantly, especially at low rates. The goal of efficient audio compression makes it imperative to optimize a correctly chosen distortion metric.

III. JOINT SELECTION OF PARAMETER VALUES: PROBLEM FORMULATION

In this section, we tackle the problem in the context of general audio coding. To concretize the presentation, we employ the AAC framework for illustrating the relevant concepts. For the general formulation, we continue to use terminology consistent with the one commonly employed in classical bit allocation, wherein the parameter values are selected for each *coefficient*. The formulation is specialized in a straightforward manner to the case of AAC, where parameter values are selected per SFB. It should perhaps be reemphasized that this approach is not restricted to AAC but is, in fact, applicable to a wide variety of audio coding standards including AC-3 [41] and G.722.1 [6].

A. Parameter Space

The quantization and encoding of each spectral coefficient is determined by a limited set of encoding parameters. In the specific case of AAC, the encoder selects values for two parameters, SF and HCB, for each SFB in the frame. Once this choice is made, the quantization and coding operations may be performed for all the coefficients in that SFB. Hence, SF and HCB, whose values are chosen per SFB, constitute the encoding parameters for AAC. The parameter space of a coefficient (or a band) is the set of *all permissible* values of all the parameters for the coefficient (or band). A point in the parameter space is given by the combination of values for the (typically multiple) encoding parameters in use by the specific compression algorithm. Note that AAC sets restrictive bounds on the quantization index values and the dynamic range of the quantized coefficients that may employ a given HCB. These restrictions effectively reduce the parameter space.

B. Cost Function Formulation

Let $p_i \in \mathbb{P}$ represent the parameter for the i th coefficient, where $\mathbb{P} = \{\psi_j, 0 < j \leq M\}$ is the parameter space with M possible parameters. Without loss of generality, we assume for simplicity the same parameter space for each coefficient.

Let the number of coefficients in the frame be N . We denote the set of parameters for the N coefficients by the vector $\mathbf{p} = (p_1, \dots, p_N) \in \mathbb{P}^N$. For the case of AAC, let us denote the set of all possible SF values by $\mathbb{S} = \{\sigma_j, 0 < j \leq M_s\}$ and HCB values by $\mathbb{H} = \{\rho_j, 0 < j \leq M_h\}$. Note that we allow for M_s distinct SF values and M_h distinct HCB values. Further, let s_i be the SF value and h_i be the HCB value for the i th SFB in the frame. Vectors \mathbf{s} and \mathbf{h} are used to denote the selected SF and HCB values for all the SFBs in the frame, i.e., $\mathbf{s} = (s_1, \dots, s_N)$ and $\mathbf{h} = (h_1, \dots, h_N)$. The combined parameter space for each SFB in AAC is the product space $\mathbb{P} = \mathbb{S} \times \mathbb{H}$ and has $M_h M_s$ elements: $\psi_j \triangleq (\sigma_j, \rho_j) \in \mathbb{S} \times \mathbb{H}, 0 < j \leq M_h M_s$.

C. Total Rate and Distortion

The total rate, R , and distortion, D , are functions of the parameter vector \mathbf{p}

$$R(\mathbf{p}) = R(p_1, \dots, p_N) \quad D(\mathbf{p}) = D(p_1, \dots, p_N). \quad (7)$$

In order to make the formulation applicable to all scenarios of potential interest, neither the distortion nor the rate is assumed additive over individual coefficients.

To illustrate this rate and distortion calculation we return to the example of AAC. The total rate required for quantization in AAC can be divided into three parts: bits required to transmit the quantized coefficient indices; bits required to transmit the SF values; and bits required to transmit the HCB values.

- Let $\mathcal{Q}(s_i, h_i)$ be the number of bits required to encode the quantized coefficient indices of the i th SFB using the SF value of s_i and HCB value of h_i . (Note that given the spectral coefficients, \mathcal{Q} is completely determined by the two parameters).
- Let \mathcal{F} denote the number of bits specifying SF for a SFB. Since AAC employs differential coding of the SFs, \mathcal{F} is a function of two parameters, s_{i-1} and s_i , for the i th SFB, and we write explicitly $\mathcal{F}(s_{i-1}, s_i)$.
- Similarly, let \mathcal{G} represent the number of bits needed to encode the HCB value of the SFB. The run-length coding of HCB produces 9 bits whenever $h_i \neq h_{i-1}$ and no bits otherwise. Hence, \mathcal{G} is a function of h_{i-1} and h_i and we write explicitly $\mathcal{G}(h_{i-1}, h_i)$.

Combining the three functions, the number of bits, b_i , for transmitting the i th SFB is given by

$$b_i(s_{i-1}, s_i, h_{i-1}, h_i) = \mathcal{Q}(s_i, h_i) + \mathcal{F}(s_{i-1}, s_i) + \mathcal{G}(h_{i-1}, h_i). \quad (8)$$

The total number of bits produced for the entire frame is then

$$\begin{aligned} R(\mathbf{s}, \mathbf{h}) &\triangleq R(\mathbf{p}), \\ &= \sum_{i=1}^N b_i(s_{i-1}, s_i, h_{i-1}, h_i), \\ &= \sum_{i=1}^N (\mathcal{Q}(s_i, h_i) + \mathcal{F}(s_{i-1}, s_i) + \mathcal{G}(h_{i-1}, h_i)) \end{aligned} \quad (9)$$

where s_0 and h_0 are initialized to zero.

Given the spectral coefficients, to calculate the distortion in SFB i we need only the band's SF value s_i , which determines the quantized coefficients, and the corresponding quantization

noise. Let $d(s_i)$ represent the quantization noise. If w_i is the weight (inverse of the masked threshold) of the i th SFB, the NMR of the SFB equals $w_i d(s_i)$. Either ANMR or MNMR can be used as the metric to combine the NMRs from the different SFBs. ANMR and MNMR for AAC can be calculated by substituting $d(s_i)$ for d_i in (1) and (2), respectively.

The problem of parameter values selection may now be stated mathematically as

$$\mathbf{p}^* = \arg \min_{\mathbf{p}: R(\mathbf{p}) \leq R_t} D(\mathbf{p}) \quad (10)$$

where R_t is the target bit rate for the frame. Note that, in the case of AAC, $R(\mathbf{p})$ is given by (9), and $D(\mathbf{p})$ by (1) or (2), depending on the criterion in use.

IV. TRELLIS-BASED OPTIMIZATION

Let us now consider the solution of the optimization problem of (10). There are M possible choices at each stage and there are N such stages. A straightforward brute-force solution to (10) has complexity in the order of $O(M^N)$. In the case of AAC, there may be as many as 49 SFBs, 60 SFs, and 12 HCBs, and the complexity of the brute-force search is $O((60 * 12)^{49})$, which is clearly impractical. We outline next an alternative approach to this optimization problem, which is based on dynamic programming [30]. First, standard Lagrangian formulation is employed to convert (10) into an unconstrained optimization problem. The Lagrangian cost function so obtained, is then demonstrated to exhibit the property of dynamic programming optimality [30]. The well-known Viterbi search [42], [43] through a trellis is applied to achieve the optimal solution at highly reduced complexity. Detailed algorithmic description of the proposed solution's application to AAC is presented for the two objective measures. (A general description of the Viterbi algorithm is available at the above references).

The standard Lagrangian procedure to reformulate the constrained optimization problem of (10) yields the Lagrangian cost

$$J(\mathbf{p}, \lambda) = D(\mathbf{p}) + \lambda R(\mathbf{p}) \quad (11)$$

where λ is the Lagrange multiplier. Clearly

$$\arg \min_{\mathbf{p}} J(\mathbf{p}, \lambda) \quad (12)$$

is the unconstrained minimization problem whose solution is also the solution of (10), once λ is adjusted to satisfy the constraint $R(\mathbf{p}) = R_t$. The original constrained minimization problem is hence solved by iterating over the different values of λ so as to achieve the target rate.

A. Dynamic Programming Solution

We construct a trellis with N stages and M states and populate the states with the parameter values $\psi_j, 0 < j \leq M$. A simple three-stage trellis is shown in Fig. 3. With every branch in this trellis we associate a cost corresponding to its contribution to the overall Lagrangian cost. The cost associated with the branch connecting p_{i-1} and p_i is denoted by $J(p_{i-1}, p_i)$. Clearly, every path through the trellis gives a particular choice of

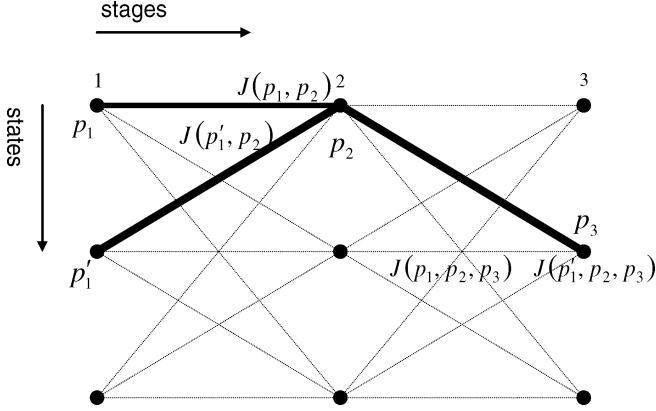


Fig. 3. Shown is a three stage, three state trellis structure in which the states represent the parameter values and the stages represent the coefficient indices. Suboptimal paths are identified and pruned employing the dynamic programming property of cost function.

encoding parameter values. We now make the standard observation: if the optimal path from p_1 to p_N passes through p_i , then it contains the optimal path from p_1 to p_i . This observation forms the basis of the Viterbi search [42], [43] and allows for an efficient search procedure where many partial paths can be pruned out without loss of optimality. The observation is shown graphically in the three-stage trellis of Fig. 3. At the second stage, $J(p_1, p_2) \leq J(p_1', p_2) \implies J(p_1, p_2, p_3) \leq J(p_1', p_2, p_3)$ and hence, $J(p_1', p_2)$ can be pruned out. Effectively, only M paths survive at the end of each stage (one ending at each state). The search then proceeds from one stage to the next and terminates at the last stage, where the entire optimal path is determined.

The use of dynamic programming leads to a dramatic reduction in complexity. Recall that the brute-force search has computational complexity of $O(M^N)$. In the dynamic programming approach, only M “best” paths are retained at any stage and comparison is carried out for N stages sequentially. For each state comparison is made from all edges branching into it (bounded by M) making the total computation complexity of the Viterbi search $O(M^2N)$, which is linear in the number of stages. Application of dynamic programming to ANMR and MNMR optimization in AAC is outlined next.

B. ANMR

The ANMR measure was discussed in Section II-A and is given by (1). The search of encoding parameter values in AAC to minimize the ANMR can be stated as

$$\arg \min_{\mathbf{s}, \mathbf{h}: R(\mathbf{s}, \mathbf{h}) \leq R_t} \frac{1}{N} \sum_{i=1}^N w_i d_i(s_i) \quad (13)$$

where $R(\mathbf{s}, \mathbf{h})$ is given by (9). The corresponding Lagrangian function is

$$\begin{aligned} J^{(A)}(\mathbf{s}, \mathbf{h}, \lambda) &= \sum_{i=1}^N w_i d_i(s_i) + \lambda \sum_{i=1}^N b_i(s_{i-1}, s_i, h_{i-1}, h_i) \\ &= \sum_{i=1}^N [w_i d_i(s_i) + \lambda (Q(s_i, h_i) + \mathcal{F}(s_{i-1}, s_i) \\ &\quad + \mathcal{G}(h_{i-1}, h_i))] \end{aligned} \quad (14)$$

where the superscript (A) indicates the ANMR measure. To summarize, the resulting optimization problem is to find the minimizer

$$(\mathbf{s}^*(\lambda), \mathbf{h}^*(\lambda)) = \arg \min_{\mathbf{s}, \mathbf{h}} J^{(A)}(\mathbf{s}, \mathbf{h}, \lambda). \quad (15)$$

We reemphasize that our problem formulation accounts for the *total number of bits* used to represent the frame, including interparameter dependencies and encoding of the side-information. Since $J^{(A)}$ is the sum of nonnegative terms and, the contribution of s_i and h_i to $J^{(A)}$ only depends on previous decisions s_{i-1} and h_{i-1} , a dynamic programming procedure can be applied to find the optimal parameter values.

The search algorithm is outlined next. A trellis is constructed where each stage corresponds to a SFB (total of 49 stages). The state j at stage i is denoted by $\Upsilon_{j,i}$. The states at a stage represent all combinations of possible choices of SF and HCB for this SFB, i.e., if the system passes through $\Upsilon_{j,i}$ then it employs the j th pair of parameter values for the i th SFB: $(s_i, h_i) = (\sigma_j, \rho_j)$. Further, we define the state-transition cost $T_{l \rightarrow j,i}$ as the cost in side-information rate for a transition from $\Upsilon_{l,i-1}$ to $\Upsilon_{j,i}$. This cost is: $T_{l \rightarrow j,i} = \lambda(\mathcal{F}(\sigma_l, \sigma_j) + \mathcal{G}(\rho_l, \rho_j))$. The minimum cost (partial) path to $\Upsilon_{j,i}$ is denoted by the vector $\Psi_{j,i}$. Finally, we denote by $J_{j,i}^{(A)}$ the cost of the minimum cost path $\Psi_{j,i}$. This is also commonly referred to as the *metric* of $\Upsilon_{j,i}$. The Viterbi search is then used to find the path through this trellis that achieves the global minimum of $J^{(A)}$ for a given λ . The value of λ that achieves the target bit rate constraint is searched using an iterative search. The search procedure is enumerated as follows.

- Step 1) *Initialize.* Set λ .
- Step 2) *Initialize.* Set metric $J_{j,0}^{(A)} = 0$, $\Psi_{j,0} = \{\emptyset\}$, $\forall j$, and $i = 1$.
- Step 3) *Search.* $\forall j$ find the best path leading to $\Upsilon_{j,i}$ by computing the metric

$$\begin{aligned} J_{j,i}^{(A)} &= \min_{l'} \left\{ J_{l,i-1}^{(A)} + w_i d_i(\sigma_j) + \lambda Q(\sigma_j, \rho_j) + T_{l \rightarrow j,i} \right\} \\ &= \min_{l'} \left\{ J_{l,i-1}^{(A)} + T_{l \rightarrow j,i} \right\} + w_i d_i(\sigma_j) + \lambda Q(\sigma_j, \rho_j) \end{aligned}$$

and let l' be the argument that achieves this minimum. The partial path leading to $\Upsilon_{j,i}$ is given by

$$\Psi_{j,i} = \{\Psi_{l',i-1}, (\sigma_j, \rho_j)\}.$$

- Step 4) *Next Stage.* If $i < N$: $i \leftarrow i + 1$, go to Step 3).
- Step 5) *Backtrack.* The best set of parameter values (overall) is given by, $(\mathbf{s}^*(\lambda), \mathbf{h}^*(\lambda)) = \Psi_{j',N}$, where $j' = \arg \min_j J_{j,N}^{(A)}$.
- Step 6) *Adjust rate.* For the optimal $\mathbf{s}^*(\lambda)$ and $\mathbf{h}^*(\lambda)$, compare total bit rate to the prescribed rate. If the constraint is not met adjust λ and go to Step 2).

C. MNMR

The MNMR measure was explained in Section II-A and is given by (2). The search of encoding parameter values in AAC to minimize the MNMR can be stated as

$$\arg \min_{\mathbf{s}, \mathbf{h}: R(\mathbf{s}, \mathbf{h}) \leq R_t} \max_{i=1}^N w_i d_i(s_i) \quad (16)$$

where $R(\mathbf{s}, \mathbf{h})$ is given by (9). The solution methodology in the MNMR case bears some similarity to that of ANMR except that, due to the min-max nature of MNMR, we do not use a classic Lagrangian approach. Instead, we define the optimal path through a trellis as the one that minimizes the rate (and purposely ignore the distortion for the moment). We hence redefine the cost function as the total rate function

$$J^{(M)}(\mathbf{s}, \mathbf{h}) = R(\mathbf{s}, \mathbf{h}) = \sum_{i=1}^N [\mathcal{Q}(s_i, h_i) + \mathcal{F}(s_{i-1}, s_i) + \mathcal{G}(h_{i-1}, h_i)]. \quad (17)$$

Since the usual observations about additivity hold for the total rate, a dynamic programming procedure can be applied to find the optimal path for the given trellis. The optimal path gives the best rate possible while ignoring the distortion incurred. The key to the solution for the MNMR case is in the construction of the trellis. Only those states are allowed (or are valid) for which the distortion is less than a certain constant value (say γ), i.e., state j in stage i is a valid state if $w_i d_i(\sigma_j) \leq \gamma$. Let $\mathbf{s}^*(\gamma)$, $\mathbf{h}^*(\gamma)$ be the set of parameter values that minimize $J^{(M)}$ when $w_i d_i(s_i) \leq \gamma, \forall i$

$$(\mathbf{s}^*(\gamma), \mathbf{h}^*(\gamma)) = \arg \min_{\mathbf{s}, \mathbf{h}} J^{(M)}(\mathbf{s}, \mathbf{h}) |_{w_i d_i(s_i) \leq \gamma}. \quad (18)$$

For such a trellis then, $\min \max_{i=1}^N w_i d_i(s_i^*) = \gamma$. Rate constraint is met by adjusting the parameter γ (not to be confused with the Lagrange multiplier λ of the ANMR case).

The search algorithm for the MNMR case is outlined next. A trellis is constructed in a fashion similar to the ANMR case, albeit with the distinction that the valid states at stage i represent all combinations of possible choices of SF and HCB values for which the NMR in the SFB is less than or equal to some constant (say γ), i.e., $\Upsilon_{j,i}$ is a valid state if $w_i d_i(\sigma_j) \leq \gamma$. Again, similar to the ANMR case, we define state-transition cost $T_{l \rightarrow j,i}$ as the cost in side-information rate for a transition from $\Upsilon_{l,i-1}$ to $\Upsilon_{j,i}$. This cost is: $T_{l \rightarrow j,i} = \mathcal{F}(\sigma_l, \sigma_j) + \mathcal{G}(\rho_l, \rho_j)$. Note the lack of the Lagrange multiplier λ in defining the state-transition cost. We also denote the minimum cost path to state $\Upsilon_{j,i}$ by the vector $\Psi_{j,i}$ and the cost of the minimum cost path by $J_{j,i}^{(M)}$. The Viterbi search is used to find the path through this trellis that achieves the global minimum of $J^{(M)}$ for a given γ . The value of γ that achieves the target bit rate constraint is searched using an iterative procedure.

- Step 1) *Initialize*. Set γ .
- Step 2) *Find Valid States*. A state $\Upsilon_{j,i}, \forall j, i$ is a valid state and retained in the trellis if $w_i d_i(\sigma_j) \leq \gamma$
- Step 3) *Initialize*. Set metric $J_{j,0}^{(M)} = 0, \Psi_{j,0} = \{\emptyset\}, \forall j$, and $i = 1$.
- Step 4) *Search*. $\forall j$ find the best path leading to $\Upsilon_{j,i}$ by computing the metric

$$J_{j,i}^{(M)} = \min_l \left\{ J_{l,i-1}^{(M)} + \mathcal{Q}(\sigma_j, \rho_j) + T_{l \rightarrow j,i} \right\}$$

and let l' be the argument that achieves this minimum. The partial path leading to $\Upsilon_{j,i}$ is given by

$$\Psi_{j,i} = \{\Psi_{l',i-1}, (\sigma_j, \rho_j)\}.$$

- Step 5) *Next Stage*. If $i < N : i \leftarrow i + 1$, go to Step 4).

Step 6) *Backtrack*. The best set of parameter values (overall) is given by, $(\mathbf{s}^*(\gamma), \mathbf{h}^*(\gamma)) = \Psi_{j',N}$, where $j' = \arg \min_j J_{j,N}^{(M)}$.

Step 7) *Adjust rate*. For the optimal $\mathbf{s}^*(\gamma)$ and $\mathbf{h}^*(\gamma)$, compare total bit rate to prescribed rate. If the constraint is not met adjust γ and go to Step 2).

For rate savings, AAC allows any set of SF and HCB values to be assigned to a SFB that is below the masking threshold. This is incorporated in our trellis by splitting every state into two—one where quantization is performed using the assigned SF and HCB values, and the other where all quantized coefficients are set to zero. The splitting of the states is similarly applied in either case of ANMR or MNMR, and results in a twofold increase in computational complexity.

V. SIMULATION RESULTS

In this section, we summarize the experimental setup including implementation details, and present the simulation results. A simplified AAC coding module derived from the publicly available MPEG AAC Verification Model (VM) [32] was employed for objective and subjective evaluation of the proposed schemes. Bit reservoir, bandwidth control and window switching modules were not employed and AAC was made to operate at a nearly constant bit rate. The implemented modules of AAC adequately serve their purpose of providing a framework for comparison of the competing search methods, albeit without attempting to achieve the performance of quality-optimized proprietary AAC encoders. For clearer comparison TLS [23] of the VM (VM-TLS) is used with some minor modification. The purpose of the modification is to emulate traditional bit allocation using the water-filling approach. VM-TLS has two nested loops. In the inner (distortion) loop, the SF values are chosen such that the NMR in each SFB is constant (say K). The total bit rate required for encoding the frame given these SFs is computed, and the value of constant K is adjusted so as to meet the constant target rate constraint in the outer (rate) loop. This modification makes no noticeable change in quality of the coded audio and no consequential difference in the rate-distortion curves of the VM-TLS presented below. The psychoacoustic model is taken from [2] and [9] with minor modifications and simplifications. The spreading function and the prediction to find the tonality factor were derived from [17] and applied to the MDCT coefficients as described in the cited reference. For the test set, eight audio files of sampling rate 44.1 kHz were taken from the EBU SQAM [32] database, which included total signals, castanets, two singing files and two speech files.

The trellis-based minimization of ANMR (TB-ANMR) and MNMR (TB-MNMR) are implemented as explained in Sections IV-B and C, respectively. The trellis states were populated with all combinations of 60 SF and 12 HCB values. Since each state was split into two, the total number of states equals $60 * 12 * 2 = 1440$. To reduce complexity, the transition at each state was restricted to the four nearest HCB values, i.e., the transition to the current state with a particular HCB value (say ρ_j) can only occur from states which have HCB values in the range $[\rho_j - 2, \rho_j + 1]$. No significant performance degradation was observed due to this restriction.

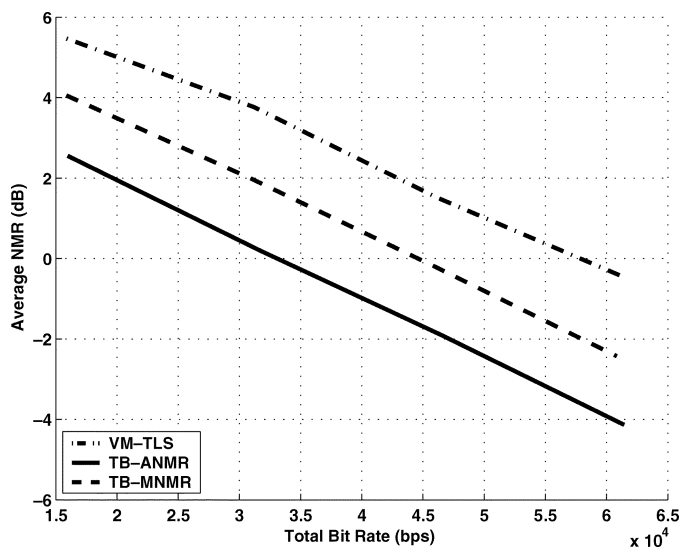


Fig. 4. Distortion-rate performance of the competing schemes. Shown is the ANMR versus bit rate for VM-TLS (dot-dashed), TB-ANMR (solid) and TB-MNMR (dashed). Note that TB-MNMR is optimized for MNMR metric but evaluated using ANMR.

A. Objective Results for a Single-Layer Coder

We compared the performance of TB-MNMR, TB-ANMR and VM-TLS on the test set. Figs. 4 and 5 depict the distortion-rate performance curves of single-layer coder over the test set. Fig. 4 shows the performance of the three schemes *evaluated using the ANMR measure*. Note specifically that TB-MNMR is optimized for the MNMR measure but evaluated here using ANMR. TB-ANMR outperforms the standard VM-TLS technique. Also of interest is the fact that the TB-MNMR scheme outperforms VM-TLS although it is evaluated using ANMR as a distortion criterion.

Fig. 5 shows the performance of three schemes evaluated using the MNMR measure. Note that, in this case, TB-ANMR is optimized for ANMR but is evaluated using MNMR. The poor performance of TB-ANMR when evaluated by the mismatched cost MNMR is explained by realizing that TB-ANMR can achieve bit rate savings by allowing high NMR in a few critical bands (and hence increase the MNMR distortion).

For both ANMR and MNMR trellis-based search outperforms the VM-TLS by a substantial margin. In particular, the performance of proposed approach yields considerable gains at coding rates of 16–48 kbps. For example, as seen from Fig. 4, TB-ANMR operating at 16 kbps achieves the same ANMR as VM-TLS at 40 kbps, while from Fig. 5, we see that TB-MNMR operating at 16 kbps achieves same MNMR as VM-TLS at 25 kbps. VM-TLS incurs a larger performance penalty when evaluated using the ANMR metric. Although VM-TLS cannot differentiate between ANMR and MNMR, simply trying to keep a constant NMR across the frequency bands results in a more severe penalty in terms of the ANMR metric than the MNMR metric, at low bit rates.

B. Subjective Results for a Single-Layer Coder

The three competing techniques were evaluated at 16 kbps using the ITU-5-grade ACR scheme [44] to produce the MOS

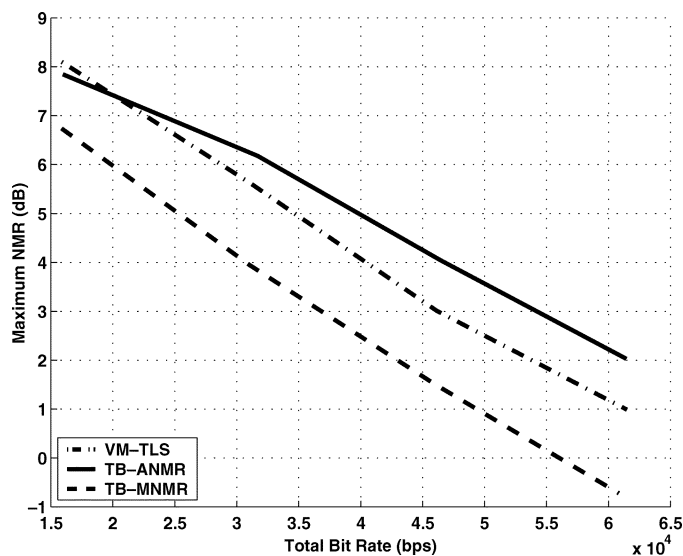


Fig. 5. Distortion-rate performance of the competing schemes. Shown is the MNMR versus bit rate for VM-TLS (dot-dashed), TB-ANMR (solid) and TB-MNMR (dashed). Note that TB-ANMR is optimized for ANMR but evaluated using MNMR.

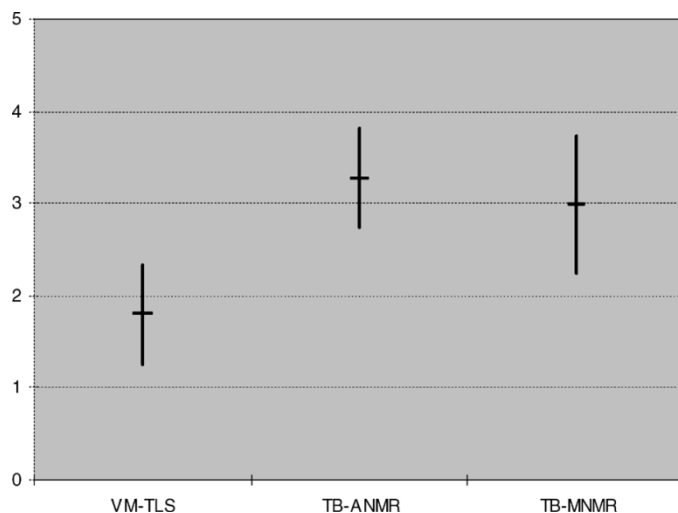


Fig. 6. Subjective five-point Mean-Opinion Score (MOS) test results for VM-TLS, TB-ANMR and TB-MNMR for a test set of eight files quantized at 16 kbps. 5 = Excellent, 4 = Good, 3 = Fair, 2 = Poor, and 1 = Bad. The vertical bars indicate 95% confidence interval. The test employed 20 listeners.

scores. The listening test was performed with 20 listeners (including several trained listeners). The test database consisted of eight files, each of about 4–8 s long. The critical test material is taken from the EBU SQAM database [31], [32] and consists of a variety of signals including German male speech, castanets, vocal singing, and harpsichord. The files were encoded using the three schemes and played twice resulting in a set of 48 occurrences (8 files \times 3 schemes \times 2 times). These 48 occurrences were played in a random order. The subjects were asked to rate each file on a 1 to 5 scale as follows: 5 = Excellent, 4 = Good, 3 = Fair, 2 = Poor, and 1 = Bad. They were also allowed to repeat the file as many times as they desired until they made their final decision. The files were played on a conventional computer with a high-end audio card using headphones. The subjective test was performed in a quiet room that was designed for audio tests. Fig. 6 shows the overall performance of the three schemes.

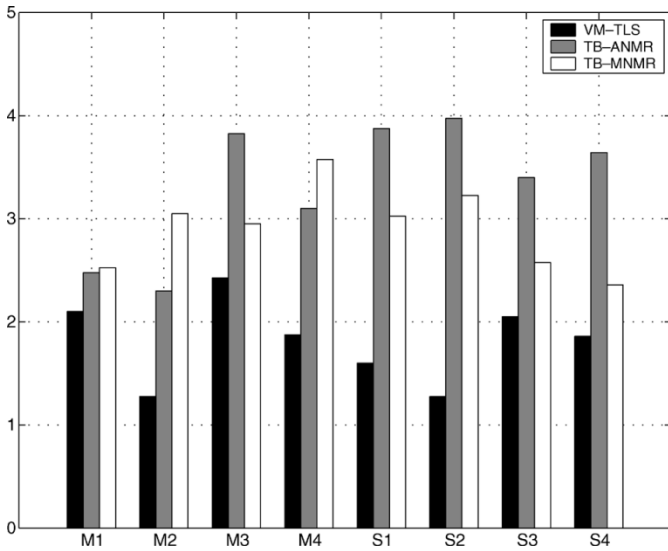


Fig. 7. Detailed break-up of subjective five-point Mean-Opinion Score (MOS) test results for VM-TLS, TB-ANMR and TB-MNMR for a test set of eight files quantized at 16 kbps. 5 = Excellent, 4 = Good, 3 = Fair, 2 = Poor, and 1 = Bad. Files M1 to M4 represent instrumental music while S1 to S4 contained vocals.

A detailed break-up of performance for each test file is given in Fig. 7. The first four files, labeled M1 to M4 are, instrumental music files and the last four, denoted by S1 to S4, are vocal singing and speech. It is clear that for vocal signals TB-ANMR performs better than TB-MNMR in all cases. For instrumental music signals, TB-MNMR performed marginally better than TB-ANMR in all cases but M3. It is interesting to note that M3 is a castanet signal containing sharp attacks. Both TB-ANMR and TB-MNMR offer substantially better quality than VM-TLS.

Furthermore, we performed an informal subjective “AB” comparison test for the TB-ANMR approach operating at 16 kbps and the VM-TLS operating at 32 kbps. The test set contained eight music and speech files. Eight listeners, some with trained ears, performed the evaluation. Each file was compressed by both competing schemes and the two compressed files were presented in random order to the listener. The listeners were asked to indicate their preference between the two samples and were also provided with the option of choosing “no preference” if no discernible difference was perceived. Within the margin of error (95% confidence interval) listeners on the average rated the overall quality of the TB-ANMR operating at 16 kbps as equivalent to that of VM-TLS operating at 32 kbps.

The nature of distortion also lends an interesting observation into the distortion metrics. As opposed to TB-ANMR, the output of TB-MNMR was mostly free of annoying artifacts such as pops and clicks. However, the output of TB-MNMR was somewhat inferior to that of TB-ANMR on the average. Most of the signals optimized using the TB-ANMR measure performed slightly better, but for a few test cases the TB-MNMR output was preferred. The resulting audio bandwidth was not substantially different in the three schemes.

C. Objective Results for a Four-Layer Scalable Coder

Fig. 8 shows the ANMR versus rate performance curves for a four-layer scalable coder where each layer operates at 16 kbps. Each layer quantizes the reconstruction error of the previous

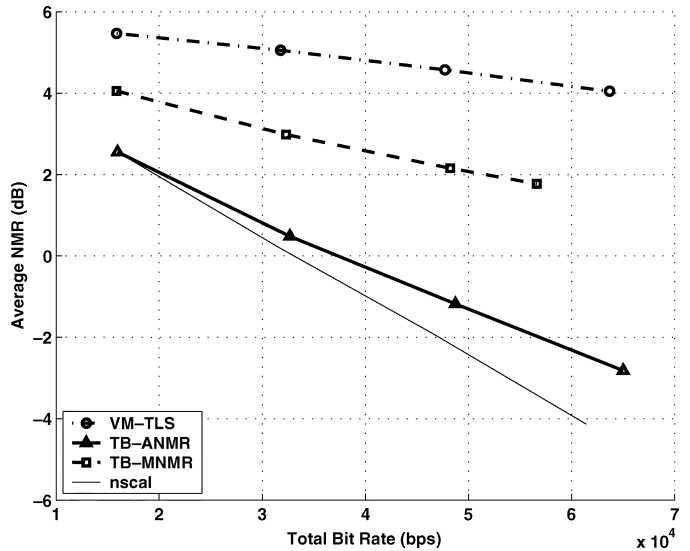


Fig. 8. Four-layer scalable coder (16/32/48/64 kbps): ANMR versus bit rate for VM-TLS and TB-ANMR and TB-MNMR. Nonscalable TB-ANMR is shown for reference.

layer. Clearly, the trellis-based approach provides major savings in bit rate over VM-TLS and these savings increase at the enhancement-layers. Also shown for reference is the nonscalable curve of TB-ANMR. This curve represents a theoretical bound on the distortion-rate performance of a scalable system. Note that the distortion-rate curve for scalable TB-ANMR approaches that of the nonscalable coder.

D. Note on Computational Complexity

The trellis-based minimization of ANMR (TB-ANMR) and MNMR (TB-MNMR) is implemented as explained in Sections IV-B and C, respectively. If the trellis states were populated with all combinations of SF and HCB values the total number of states equals 1440 and the complexity of the full trellis search is in the order of two million operations per SFB. (The trellis complexity is linear rather than exponential in the number of stages or SFBs, but quadratic in the number of states.) Recall that to reduce the complexity further in our simulations, the transition at each state was restricted to the four nearest HCB values. Hence, the search complexity for the trellis-based scheme is reduced by another factor of three. In two recent papers [33], [34], the authors propose and discuss approaches for further reduction of the search complexity. A nonoptimized implementation of the proposed trellis-based scheme on a Pentium 1.6-GHz machine was 25 times more complex as compared to VM-TLS approach.

VI. CONCLUSION

In this paper, we derived a trellis-based optimization scheme for AAC for minimizing two different objective measures; average NMR and maximum NMR. The scheme substantially enhances performance at low bit rates. Under parameter independence and high-resolution assumptions, the two objective measures yield an identical solution. However, ignoring parameter dependencies leads to poor performance at low rates. The main contributions were the reformulation of the parameter optimization problem at the encoder to account for interparameter dependencies in encoding side-information, and the development of a dynamic programming technique to

obtain the solution at manageable complexity. The resulting bit stream is standard-compatible, and the additional computational complexity is incurred only at the encoder. Simulation results employing AAC on the SQAM database demonstrate considerable gains at low bit rates.

APPENDIX

This Appendix sketches briefly a demonstration that, under the assumptions of high-resolution and interband parameter independence, ANMR and MNMR lead to the same solution. The bit allocation problem for minimizing ANMR is defined as

$$\min_{\mathbf{b}} \sum_{i=1}^N d_i(b_i) \quad \text{such that} \quad \sum_{i=1}^N b_i \leq R_t. \quad (19)$$

The bit allocation problem for minimizing MNMR is defined as

$$\min_{\mathbf{b}} \max_{i=1}^N d_i(b_i) \quad \text{such that} \quad \sum_{i=1}^N b_i \leq R_t. \quad (20)$$

The high-resolution model-based solution for the ANMR measure is obtained by the minimizer given in (5). We claim that using the high-resolution distortion model of (4), the solution to the MNMR problem of (20) results in the same minimizer of (5), which we repeat here

$$b_i^* = \frac{R_t}{N} + \frac{1}{2} \log_2 \frac{\nu_i^2}{\omega^2}.$$

The claim is proven based on a simple argument: Let \mathbf{b}^* be the ANMR minimizer. By the water filling principle, or equal distortion in all bands as given in (6), we may write explicitly that $d_i(b_i^*) = d$, a constant for all i . Next, let \mathbf{b} be any other assignment such that $\sum_{i=1}^N b_i \leq R_t$ (a tentative MNMR solution). By (19), we know that

$$\frac{1}{N} \sum_{i=1}^N d_i(b_i) \geq \frac{1}{N} \sum_{i=1}^N d_i(b_i^*) = d.$$

This implies immediately that

$$\max_{i=1}^N d_i(b_i) \geq d = \max_{i=1}^N d_i(b_i^*) \quad (21)$$

i.e., \mathbf{b}^* is also the MNMR solution.

REFERENCES

- [1] *Information Technology—Generic Coding of Moving Pictures and Associated Audio*, ISO/IEC Std. ISO/IEC JTC1/SC29 13 818-7:1997(E), 1997.
- [2] *Information Technology—Very Low Bitrate Audio-Visual Coding*, ISO/IEC Std. ISO/IEC JTC1/SC29 14 496-3:2001(E), 2001.
- [3] L. D. Fielder, M. Bosi, G. Davidson, M. Davis, C. Todd, and S. Vernon, "AC-2 and AC-3: low-complexity transform-based audio coding," in *Collected Papers on Digital Audio Bit-Rate Reduction*, N. Gilchrist and C. Grewin, Eds. New York: Audio Eng. Soc., 1996, pp. 54–72.
- [4] D. Sinha, J. D. Johnston, S. Dorward, and S. Quackenbush, "The perceptual audio coder (PAC)," in *Digital Signal Processing Handbook*, V. Madiseti and D. B. Williams, Eds. New York: IEEE Press, 1998.
- [5] K. Akagiri, M. Katakura, H. Yamauchi, E. Saito, M. Kohut, M. Nishiguchi, and K. Tsutsui, "Sony systems," in *Digital Signal Processing Handbook*, V. Madiseti and D. B. Williams, Eds. New York: IEEE Press, 1998.
- [6] *Coding at 24 and 32 kbit/s for Hands-Free Operation in Systems With Low Frame Loss*, ITU-T Std. ITU-T Recommendation G.722.1, Sep. 1999.
- [7] E. Zwicker and H. Fastl, *Psychoacoustics: Facts and Models*, 2nd ed. New York: Springer-Verlag, 1999.
- [8] H. Fletcher, "Auditory patterns," *Rev. Modern Phys.*, vol. 12, pp. 47–65, Jan. 1940.
- [9] J. D. Johnston, "Transform coding of audio signals using perceptual noise criteria," *IEEE J. Select. Areas Commun.*, vol. 6, no. 2, pp. 314–323, Feb. 1988.
- [10] M. R. Schroeder, B. S. Atal, and J. L. Hall, "Optimizing digital speech coders by exploiting masking properties of the human ear," *J. Acoust. Soc. Amer.*, vol. 66, no. 6, pp. 1647–1652, Dec. 1979.
- [11] T. Painter and A. Spanias, "Perceptual coding of digital audio," *Proc. IEEE*, vol. 88, no. 4, pp. 451–515, Apr. 2000.
- [12] K. Ramchandran, A. Ortega, and M. Vetterli, "Bit allocation for dependent quantization with applications to multiresolution and MPEG video coders," *IEEE Trans. Image Process.*, vol. 3, no. 5, pp. 533–545, Sep. 1994.
- [13] Y. Shoham and A. Gersho, "Efficient bit allocation for an arbitrary set of quantizers," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 36, no. 9, pp. 1445–1453, Sep. 1988.
- [14] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Norwell, MA: Kluwer, 1992.
- [15] R. J. Beaton, J. G. Beerends, M. Keyhl, and W. C. Treurniet, "Objective perceptual measurement of audio quality," in *Collected Papers on Digital Audio Bit-Rate Reduction*, N. Gilchrist and C. Grewin, Eds. New York: Audio Eng. Soc., 1996, pp. 126–152.
- [16] K. Brandenburg, "Evaluation of quality for audio encoding at low bit rates," in *Proc. 82nd AES Convention*, 1987.
- [17] H. Najafzadeh-Azghandi and P. Kabal, "Improving perceptual coding of narrow-band audio signals at low rates," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 2, Mar. 1999, pp. 913–916.
- [18] H. Najafzadeh and P. Kabal, "Perceptual bit allocation for low rate coding of narrow-band audio," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 2, Jun. 2000, pp. 893–896.
- [19] N. S. Jayant and P. Noll, *Digital Coding of Waveforms: Principles and Applications to Speech and Video*. Englewood Cliffs, NJ: Prentice-Hall, 1984.
- [20] J. J. Huang and P. M. Schultheiss, "Block quantization of correlated Gaussian random variables," *IEEE Trans. Commun. Syst.*, vol. CS-1, pp. 289–296, Sep. 1963.
- [21] B. Fox, "Discrete optimization via marginal analysis," *Manage. Sci.*, vol. 13, no. 3, pp. 210–216, Nov. 1966.
- [22] A. Segall, "Bit allocation and encoding for vector sources," *IEEE Trans. Inf. Theory*, vol. IT-22, no. 2, pp. 162–169, Mar. 1976.
- [23] M. Bosi, K. Brandenburg, S. Quackenbush, L. Fielder, K. Akagiri, H. Fuchs, M. Dietz, J. Herre, G. Davidson, and Y. Oikawa, "ISO/IEC MPEG-2 advanced audio coding," *J. Audio Eng. Soc.*, vol. 45, no. 10, pp. 789–814, Oct. 1997.
- [24] P. A. Chou, T. Lookabaugh, and R. M. Gray, "Optimal pruning with applications to tree-structured source coding and modeling," *IEEE Trans. Inf. Theory*, vol. 35, no. 2, pp. 299–315, Mar. 1989.
- [25] E. A. Riskin, "Optimal bit allocation via the generalized BFOS algorithm," *IEEE Trans. Inf. Theory*, vol. 37, no. 2, pp. 400–402, Mar. 1991.
- [26] P. Prandoni and M. Vetterli, "Optimal bit allocation with side information," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 5, Jun. 1999, pp. 2411–2414.
- [27] L. P. Kondi and A. K. Katsaggelos, "An operational rate-distortion optimal single-pass SNR scalable video coder," in *Proc. IEEE Int. Conf. Image Processing*, vol. 10, Nov. 2001, pp. 1613–1620.
- [28] B. Grill, "A bit rate scalable perceptual coder for MPEG-4 audio," in *Proc. 103rd AES Convention*, New York, 1997.
- [29] —, "Scalable joint stereo coding," in *Proc. 105th AES Convention*, San Francisco, CA, 1998.
- [30] R. E. Bellman, *Dynamic Programming*. Princeton, NJ: Princeton Univ. Press, 1957.
- [31] *Sound Quality Assessment Material Recordings for Subjective Tests*, European Broadcasting Union (EBU) Std., Rev. Tech. 3253-E, Apr. 1988.
- [32] The MPEG Audio Web Page. [Online]. Available: <http://www.tnt.uni-hannover.de/project/mpeg/audio/>
- [33] C.-H. Yang and H.-M. Hang, "Efficient bit assignment strategy for perceptual audio coding," in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 5, Apr. 2003, pp. 405–408.
- [34] —, "Cascaded trellis-based optimization for MPEG-4 advanced audio coding," in *Proc. 115th AES Convention*, New York, 2003.
- [35] A. Aggarwal, S. L. Regunathan, and K. Rose, "Trellis-based optimization of MPEG-4 advanced audio coding," in *Proc. IEEE Workshop on Speech Coding*, Sep. 2000, pp. 142–144.

- [36] —, “Near-optimal selection of encoding parameters for audio coding,” in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 5, May 2001, pp. 3269–3272.
- [37] J. G. Beerends and J. A. Stemerdink, “A perceptual audio quality measure based on a psychoacoustic sound representation,” *J. Audio Eng. Soc.*, vol. 40, no. 12, pp. 963–978, Dec. 1992.
- [38] W. C. Treurniet and G. A. Souloudre, “Evaluation of the ITU-R objective audio quality measurement method,” *J. Audio Eng. Soc.*, vol. 48, no. 3, pp. 164–173, Mar. 2000.
- [39] *Method for Objective Measurements of Perceived Audio Quality*, ITU-R Std. BS.1387-1, Nov. 2001.
- [40] Y. Shoham and A. Gersho, “Efficient codebook allocation for an arbitrary set of vector quantizers,” in *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. 4, 1985, pp. 1696–1699.
- [41] *Digital Audio Compression Standard (AC-3)*, ATSC Std. A/52, Dec. 1995.
- [42] A. J. Viterbi, “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm,” *IEEE Trans. Inf. Theory*, vol. IT-13, no. 4, pp. 260–269, Apr. 1967.
- [43] G. D. Forney, Jr., “The Viterbi algorithm,” *Proc. IEEE*, vol. 61, no. 3, pp. 268–278, Mar. 1973.
- [44] *Methods for Subjective Determination of Transmission Quality*, ITU-T Std., Rev. Recommend. P.800, Aug. 1996.



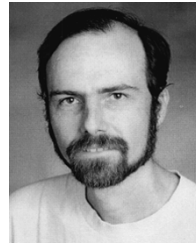
Ashish Aggarwal (S'97–M'93) received the B.E degree in electronics from Bombay University, Bombay, India, in 1996, and the M.S. and Ph.D. degrees in electrical engineering from the University of California, Santa Barbara, in 1998 and 2002, respectively.

From July 2002 to July 2003, he was employed by PortalPlayer, Inc., where he carried out the design and implementation of audio coders such as MP3 and AAC. In July 2003, he joined Harman International's Advanced Technology Group. His main research activities are audio compression and post processing algorithms.

Dr. Aggarwal currently serves as a member of the IEEE Technical Committee on Audio and Electroacoustics. He is a member of the Signal Processing and Communications Societies of the IEEE and a member of the AES.

Shankar L. Regunathan (S'96–M'01) received the B.Tech degree in electronics and Communication from the Indian Institute of Technology, Madras, in 1994, and the M.S. and Ph.D. degrees in electrical engineering from the University of California, Santa Barbara, in 1996 and 2001, respectively.

Currently, he is with Microsoft Corporation, Redmond, WA.



Kenneth Rose (S'85–M'91–SM'01–F'03) received the B.Sc. degree (summa cum laude) and M.Sc. degree (magna cum laude) in electrical engineering from Tel-Aviv University, Tel-Aviv, Israel, in 1983 and 1987, respectively. In 1991 he received the Ph.D. degree from the California Institute of Technology, Pasadena.

From July 1983 to July 1988, he was employed by Tadiran, Ltd., Israel, where he carried out research in the areas of image coding, image transmission through noisy channels, and general image processing. In January 1991, he joined the Department of Electrical and Computer Engineering, University of California at Santa Barbara, where he is currently a Professor. His main research activities are in information theory, source and channel coding, video and audio coding and networking, pattern recognition, and nonconvex optimization in general. He is also particularly interested in the relations between information theory, estimation theory, and statistical physics, and their potential impact on fundamental and practical problems in diverse disciplines.

Dr. Rose currently serves as Area Editor for the IEEE TRANSACTIONS ON COMMUNICATIONS. He cochaired the technical program committee of the 2001 IEEE Workshop on Multimedia Signal Processing. In 1990 he received (with A. Heiman) the William R. Bennett Prize Paper Award from the IEEE Communications Society.