

Orthonormal and Biorthonormal Filter Banks as Convolvers, and Convolutional Coding Gain

P. P. Vaidyanathan, *Fellow, IEEE*

Abstract—A maximally decimated filter bank system (with possibly unequal decimation ratios in the subbands) can be regarded as a generalization of the short-time Fourier transformer. In fact, it is known that such a “filter bank transformer” is closely related to the wavelet transformation. A natural question that arises when we conceptually pass from the traditional Fourier transformer to the filter bank transformer is: what happens to the convolution theorem, i.e., is there an analog of the convolution theorem in the world of “filter bank transforms”? In this paper we address the question first for uniform decimation and then generalize it to the nonuniform case. The result takes a particularly simple and useful form for paraunitary or orthonormal filter banks. It shows how we can convolve two signals $x(n)$ and $g(n)$ by directly convolving the subband signals of a paraunitary filter bank and adding the results. The advantage of the method is that we can quantize in the subbands based on the signal variance and other perceptual considerations, as in traditional subband coding. As a result, for a fixed bit rate, the result of convolution is much more accurate than direct convolution. That is, we obtain a coding gain over direct convolution. We will derive expressions for optimal bit allocation and optimal coding gain for such paraunitary convolvers. As a special case, if we take one of the two signals to be the delta function (e.g., $g(n) = \delta(n)$), we can recover the well-known bit allocation and coding gain formulas of traditional subband coding. The derivations also show that these formulas are valid regardless of the filter quality, as long as orthonormality is not violated.

A special case similar to orthogonal transform coding is also considered and good convolutional coding gains for speech are demonstrated, with the use of the DCT matrix. Finally, the result is extended to the case of nonuniform biorthonormal filter banks; with the incorporation of an additional trick, the convolution theorems in this case become as simple as for the orthonormal case.

I. INTRODUCTION

SHOWN in Fig. 1(a) is the M -channel maximally decimated digital filter bank, which has been studied by a number of authors in the past decade. Here $H_k(z)$, $F_k(z)$, $0 \leq k \leq M-1$ are the set of analysis and synthesis filters. The notations \downarrow_{n_k} and \uparrow_{n_k} denote the n_k -fold decimator and interpolator (unpsampler or expander) as defined in several earlier references [1]–[5]. In this paper,

Manuscript received March 11, 1992; revised July 26, 1992. The associate editor coordinating the review of this paper and approving it for publication was Dr. S. D. Cabrera. This work was supported in part by National Science Foundation Grant MIP 8919196 with matching funds from Rockwell Inc. and Tektronix, Inc.

The author is with the Department of Electrical Engineering, California Institute of Technology, Pasadena, CA 91125.

IEEE Log Number 9208200.

all n_k are positive integers. The boxes labeled Q_k denote quantizers which quantize the subband signals $x_k(n)$.

The relations between filter banks and wavelet transforms have been known for some time [6]–[12]. An excellent magazine article appeared recently [10], discussing this connection explicitly. It is well known that wavelet transforms provide more flexibility (in terms of time-frequency resolution) than the traditional Fourier transform. In this paper we deal only with discrete-time filter banks (both uniform and nonuniform decimators will be considered). It is known that discrete time filter banks can be considered as discrete time wavelet transformations. Here the analysis bank can be viewed as a transformation from “time” to “time-frequency.” We will simply refer to this as the filter bank transform, and the decimated subband signals $x_k(n)$ will be called the transform-domain signals. The synthesis bank is regarded as the inverse transformer (assuming perfect reconstruction, that is, $\hat{x}(n) = x(n)$).

A. Aim of the Paper

The advent of these transforms leads us to ask the question, how do the standard properties of the Fourier transformation generalize to the case of “filter-bank transforms”? For example, what is the extension of the convolution theorem? To introduce the main topic of the paper, let $y(n)$ denote the convolution of two sequences $x(n)$ and $g(n)$, i.e., $y(n) = \sum_{m=-\infty}^{\infty} x(m)g(n-m)$. According to Fourier transform theory [13], the transform of $y(n)$ is related to those of $x(n)$ and $g(n)$ as $Y(e^{j\omega}) = X(e^{j\omega})G(e^{j\omega})$, i.e., convolution becomes “multiplication” in the transform domain. Now consider the “filter bank transformer,” with the decimated subband signals regarded as the “transform domain.” What is the “convolution theorem” in this case? To expand on this question, consider Fig. 1 where we show $x(n)$ and $g(n)$ as the inputs to two copies of the filter bank. The transform domain “coefficients” corresponding to $x(n)$ and $g(n)$ are the sets of sequences $x_k(n)$ and $g_k(n)$, respectively. How should we combine $x_k(n)$, $g_k(n)$, $0 \leq k \leq M-1$ so that the convolution $\sum_m x(m)g(n-m)$ can be obtained from these, assuming there are no subband quantizers?

In Section II-A we will derive this convolution theorem for the case of uniform filter banks (i.e., $n_k = M$ for all k). The result takes an exceptionally simple form in the case of paraunitary or orthonormal filter banks [2], [12],

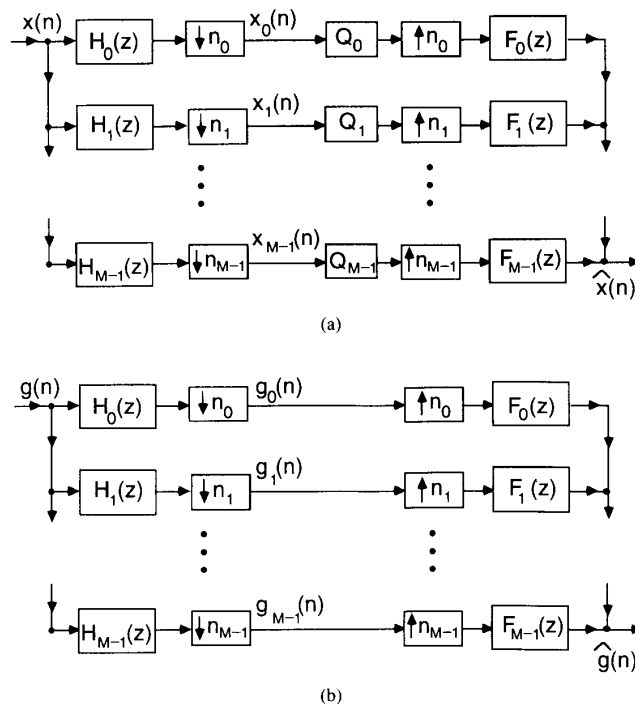


Fig. 1. The maximally decimated filter bank: (a) with input $x(n)$, and (b) with input $g(n)$.

[14], [15]. Qualitatively speaking, the convolution $x(n) * g(n)$ is reduced to computing the convolutions $x_k(n) * g_k(n)$ and adding. In other words, the subband convolutions are decoupled; we need not consider $x_k(n) * g_m(n)$, for $k \neq m$. The result will be stated more precisely in Theorem 2.1 (equal n_k) and Theorem 2.2 (unequal n_k). A similar result also follows for the nonorthonormal case (the biorthonormal case), with the incorporation of a simple additional trick, as shown in Section VI (Theorem 6.1).

In Section II-B, the result will be extended to the case of filter banks with nonuniform decimation ratios. Once again, it will be shown that when the synthesis filter coefficients form an orthonormal basis (this being the extension of the paraunitary concept), the "convolution theorem" takes a special simple form. Even though the uniform filter bank is a special case, we have chosen to treat it separately first, because it is much simpler, while conveying most of the ideas well.

B. Usefulness

The motivation for obtaining these "convolution theorems" does not originate from a desire to obtain algorithms that are faster than the many well-known fast convolution techniques. (Indeed, the state of the art for fast convolutions is already very advanced.) The actual motivation comes from the fact that we can quantize in the subbands, and reduce the roundoff error (for fixed wordlength) by the proper bit allocation schemes. Thus,

instead of quantizing $x(n)$ and then convolving with $g(n)$, we can now quantize $x_k(n)$ and then convolve with $g_k(n)$ and add the results for all k . When performing this quantization in subbands, we can exploit the subband energy distribution and perform optimal bit allocation. In this way, we obtain increased accuracy for a given bit rate. That is, the system offers a coding gain. This idea is very similar in philosophy to subband coding [16] (e.g., see [17, ch. 11] and [18, ch. 1]).

Unlike subband coding where subband quantization is used to compress the amount of data to be transmitted or stored, the goal is somewhat different here. The quantization of subband signals here allows the subband convolutions to be implemented faster, with certain types of computational architectures (e.g., bit-serial). Clearly the usefulness in a particular application depends on the chosen architecture, demands on speed and accuracy, and so forth.

In a spirit to that described in the above references, we can define a coding gain for the paraunitary convolver. We will obtain the optimal bit allocation formula, and study the coding gain under optimal bit allocation. Unlike in usual subband coding, it is possible to obtain a coding gain > 1 even if $x(n)$ has a flat spectrum (i.e., is white).

It is important to notice that the computation of the subband signals $x_k(n)$ itself involves filtering. If this filtering complexity is comparable to the direct convolution of $x(n)$ and $g(n)$, then the above technique is clearly unworthy. It has potential applications when $x(n)$ and $g(n)$ are very

long sequences (in comparison with the lengths of the analysis filters $H_k(z)$). A useful special case arises when the analysis filters have length $\leq M$ (which is analogous to transform coding, e.g., using the DCT). We will see in Section IV that, even in this case, substantial coding gain can be exhibited.

C. Paper's Outline and Main Results

1) *Orthonormal case.* In Section II-A we derive the convolution theorem for paraunitary filter banks with uniform decimation. This is extended to the case of nonuniform filter banks (n_k not identical for all k) in Section II-B. We show how orthonormality can be exploited to decouple the subband convolutions.

2) *Bit allocation and convolutional coding gain.* Section III presents a derivation of optimal subband bit allocation, as well as the corresponding coding gain expression for the orthonormal convolver. Both uniform and nonuniform cases will be considered.

3) *Traditional subband coding.* In Section III-D we show how the well-known coding gain results for traditional subband systems can be obtained as special cases of the convolver's optimal bit allocation and coding gain expressions. (For this, one of the two signals to be convolved is taken as the unit pulse.) This also shows that the ideal brick-wall nature of the analysis filters is not necessary, as is sometimes assumed, for the validity of these expressions; paraunitariness (more generally orthonormality) is sufficient.

4) *Transform coding, and extension of the KLT problem.* Section IV considers a different specialization of the uniform paraunitary convolver, with analysis filter lengths constrained to be $\leq M$. This is, in principle, the extension of the transform coding problem to the case of convolution. It has the advantage that we can further maximize the coding gain of the optimally bit-allocated system, by optimizing the transform matrix (generalization of the KLT). Section V presents several numerical examples.

5) *Biorthonormal case.* In Section VI we show how the convolution theorem can be extended to the filter banks that do not satisfy orthonormality, but only the perfect reconstruction (or the biorthonormality) property. Even in this case we show that the subband convolutions can be decoupled, provided we use the analysis filters to decompose $x(n)$, and the synthesis filters to decompose $g(n)$ (see Fig. 11 for a preview).

D. Notations and Basics

1) $x(n) * g(n)$ denotes convolution of $x(n)$ with $g(n)$. The sequence $x(n) * g^*(-n)$ is the deterministic cross correlation between $x(n)$ and $g(n)$, and has z -transform $X(z)\bar{G}(z)$. All signals and impulse responses are, in general, infinitely long and possibly noncausal.

2) Boldfaced quantities represent matrices and vectors. The notations A^T , A^* , and A^\dagger represent, respectively, the transpose, conjugate, and transpose-conjugate of A . The tilde, as in $\tilde{H}(z)$, stands for transposition, followed by

conjugation of coefficients, followed by replacement of z with z^{-1} . Thus $H(z) = \sum_n h(n)z^{-n}$ implies $\tilde{H}(z) = \sum_n h^\dagger(-n)z^{-n}$. For any z , we have $\tilde{H}(z) = H^\dagger(1/z^*)$; on the unit circle $\tilde{H}(z) = H^\dagger(z)$.

3) $W_N = e^{-j2\pi/N}$, with subscript omitted when it is clear.

4) The M -fold decimator $\downarrow M$ and expander $\uparrow M$ (or interpolator) are defined as in [1], [2]. Thus the input output relation for the decimator is $y(n) = x(Mn)$, and for the expander it is

$$y(n) = \begin{cases} x(n/M), & n = \text{integer mul. of } M \\ 0, & \text{otherwise.} \end{cases}$$

In this paper, all decimation and interpolation ratios are positive integers. In equations, the notation $a(n)|_{\downarrow M}$ denotes the decimated sequence $a(Mn)$. (The vertical bar is omitted where it is unnecessary.) With $A(z)$ denoting the z transform of $a(n)$, the notation $A(z)|_{\downarrow M}$ denotes the z transform of the decimated version $a(Mn)$. Let $A(z)$ and $B(z)$ be rational functions and let K and L be integers. The following identity can be easily verified:

$$(A(z^K)B(z))_{\downarrow KL} = (A(z)(B(z)|_{\downarrow K}))_{\downarrow L}. \quad (1.1)$$

E. Polyphase Notation

For the case where $n_k = M$ for all k , the system of Fig. 1(a) can be redrawn as in Fig. 2 where $E(z)$ and $R(z)$ are $M \times M$ matrices. Defining the analysis and synthesis filter vectors as

$$\begin{aligned} \mathbf{h}(z) &= [H_0(z) \ H_1(z) \ \cdots \ H_{M-1}(z)]^T, \\ \mathbf{f}(z) &= [F_0(z) \ F_1(z) \ \cdots \ F_{M-1}(z)]^T \end{aligned} \quad (1.2)$$

we have

$$\mathbf{h}(z) = E(z^M)\mathbf{e}(z), \quad \mathbf{f}^T(z) = \tilde{\mathbf{e}}(z)\mathbf{R}(z^M) \quad (1.3)$$

where $\mathbf{e}(z)$ is the delay chain vector, i.e.,

$$\mathbf{e}(z) = [1 \ z^{-1} \ \cdots \ z^{-(M-1)}]^T. \quad (1.4)$$

$E(z)$ and $R(z)$ are, respectively, the polyphase matrices of the analysis and synthesis banks. Note, in particular, that any transfer function $H(z)$ can be written in the form

$$H(z) = \sum_{n=0}^{M-1} z^{-n} E_n(z^M)$$

where $E_n(z)$ are the so-called Type 1 polyphase components.

II. CONVOLUTION THEOREMS FOR ORTHONORMAL FILTER BANKS

A. Filter Bank with Equal Decimation Ratio in all Branches

First consider Fig. 1 with $n_k = M$ for all k . The convolution theorem is obtained by analyzing this in absence of the quantizers Q_k . Assume that the set of filters $\{H_k(z), F_k(z)\}$ are chosen to satisfy the perfect reconstruction

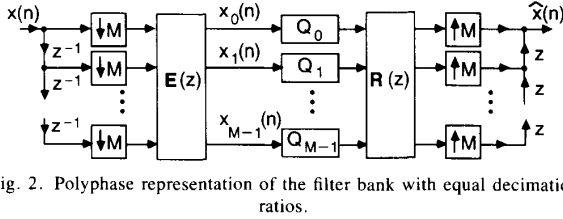


Fig. 2. Polyphase representation of the filter bank with equal decimation ratios.

property, i.e.,

$$\hat{X}(z) = X(z), \quad \hat{G}(z) = G(z). \quad (2.1)$$

Using the fact that the M -fold upsamplers have outputs $X_k(z^M)$ and $G_k(z^M)$, we can express $\hat{X}(z)$ as $\sum_{k=0}^{M-1} X_k(z^M) F_k(z)$, and similarly for $\hat{G}(z)$. Using these together with (2.1) we obtain

$$\begin{aligned} X(z) &= \sum_{k=0}^{M-1} X_k(z^M) F_k(z) \\ G(z) &= \sum_{k=0}^{M-1} G_k(z^M) F_k(z). \end{aligned} \quad (2.2)$$

Now consider the quantity $X(z) \tilde{G}(z)$ (with the tilde as defined at the end of Section I). We have

$$X(z) \tilde{G}(z) = \sum_{k=0}^{M-1} \sum_{m=0}^{M-1} X_k(z^M) \tilde{G}_m(z^M) F_k(z) \tilde{F}_m(z). \quad (2.3)$$

The inverse z transform of $X(z) \tilde{G}(z)$ is equal to the convolution of $x(n)$ with $g^*(-n)$ (i.e., the deterministic cross correlation between $x(n)$ and $g(n)$). Similarly $X_k(z) \tilde{G}_m(z)$ represents the convolution of the subband signals $x_k(n)$ and $g_m^*(-n)$.

1) *Paraunitary or Orthonormal Filter Banks:* The above equation reduces to a much simpler form (the double summation reduces to a single summation) when the filter bank is paraunitary [12], [14], [15].¹ In this case the polyphase matrix $E(z)$ satisfies

$$\tilde{E}(z) E(z) = I \quad (2.4)$$

and we choose $R(z) = \tilde{E}(z)$ for perfect reconstruction (so that $R(z)$ is also paraunitary). In this case the analysis and synthesis filters are related as $F_k(z) = \tilde{H}_k(z)$, that is, $f_k(n) = h_k^*(-n)$. In order to ensure that $F_k(z)$ is stable, we assume that the analysis filters are FIR. Thus, $h_k(n)$ and $f_k(n)$ are FIR with the same length. A paraunitary filter bank satisfies the following properties, regardless of the exact nature of $H_k(e^{j\omega})$ (i.e., regardless of filter quality) [12].

1) The energy of each analysis filter equals unity, that is $\int_0^{2\pi} |H_k(e^{j\omega})|^2 d\omega / 2\pi = 1$.

2) The analysis filters are power complementary, that is, $\sum_k |H_k(e^{j\omega})|^2 = M$.

3) Since $f_k(n) = h_k^*(-n)$, we have $|F_k(e^{j\omega})| = |H_k(e^{j\omega})|$. So the above two properties hold for the synthesis filters as well.

¹To appreciate the significance of the simplification, read Section VII-A.

(Notice, in particular, that in the case of ideal brick-wall filters, to be shown later in Fig. 6, the first two properties are evident.) The paraunitary property of $R(z)$ is equivalent to the property that the synthesis filters satisfy an orthonormality condition [10]–[12], that is,

$$\sum_{n=-\infty}^{\infty} f_k(n) f_m^*(n + Mi) = \delta(k - m) \delta(i). \quad (2.5)$$

In the z -domain this can be rewritten as

$$(F_k(z) \tilde{F}_m(z))_{\downarrow M} = \delta(k - m) \quad (\text{orthonormality}). \quad (2.6)$$

2) *Simplification of the Convolution Formula:* Using the above orthonormality condition, (2.3) leads to

$$(X(z) \tilde{G}(z))_{\downarrow M} = \sum_{k=0}^{M-1} X_k(z) \tilde{G}_k(z). \quad (2.7)$$

This can be rewritten in the time domain as

$$(x(n) * g^*(-n))_{\downarrow M} = \sum_{k=0}^{M-1} x_k(n) * g_k^*(-n). \quad (2.8)$$

In the time domain, the left-hand side represents the M -fold decimated version of the convolution of $x(n)$ with $g^*(-n)$. The k th term on the right-hand side represents the convolution of the subband signal $x_k(n)$ with $g_k^*(-n)$. The cross terms of the form $x_k(n) * g_m^*(-n)$ have disappeared, i.e., there is no “cross-coupling” between subbands any more. Summarizing, we have proved:

*Theorem 2.1: Paraunitary convolution theorem.*² Consider the two copies of a maximally decimated filter bank as in Fig. 1, with FIR analysis and synthesis filters, and $n_k = M$ for all k . Ignore the quantizers Q_k . Assume that the system has perfect reconstruction ($\hat{x}(n) = x(n)$ for any $x(n)$) and that the polyphase matrix $E(z)$ (Fig. 2) is paraunitary (equivalently the synthesis filters are orthonormal, i.e., satisfy (2.5) or equivalently (2.6)). Then the M -fold decimated version of the convolution $x(n) * g^*(-n)$ can be computed by computing the convolutions $x_k(n) * g_k^*(-n)$, $0 \leq k \leq M-1$, and adding them. \diamond

Obtaining all the samples: In order to obtain all the samples of the convolution $x(n) * g^*(-n)$, we repeat (at least conceptually) the above operation M times, by replacing $g(n)$ with $g(n - i)$, for $0 \leq i \leq M-1$. We can represent these operations mathematically as

$$(z^i X(z) \tilde{G}(z))_{\downarrow M} = \sum_{k=0}^{M-1} X_k(z) \tilde{G}_k^{(i)}(z), \quad 0 \leq i \leq M-1 \quad (2.9)$$

where $G_k^{(i)}(z)$ is the subband signal obtained by replacing $g(n)$ with $g(n - i)$. Equation (2.9) means that the i th Type 1 polyphase component (see end of Section I) of the quantity $X(z) \tilde{G}(z)$ is given by the right-hand side. So we can write

$$X(z) \tilde{G}(z) = \sum_{k=0}^{M-1} X_k(z^M) \sum_{m=0}^{M-1} z^{-m} \tilde{G}_k^{(m)}(z^M). \quad (2.10)$$

²A simple modification of this result, which eliminates the need for paraunitariness and depends only on the so-called biorthonormality, is presented in Section VI. The modification also shows how we can directly perform convolution $(x(n) * g(n))$ rather than correlation $(x(n) * g^*(-n))$.

Notice that it is not necessary to repeatedly run the filter bank with inputs $g(n-i)$ for $0 \leq i \leq M-1$ in order to obtain $g_k^{(i)}(n)$. Let $s_k(n)$ denote the (undecimated) output of $H_k(z)$ in response to $g(n)$. Then $g_k^{(i)}(n) = s_k(Mn-i)$ so that we can write $S_k(z) = \sum_{m=0}^{M-1} z^m G_k^{(m)}(z^M)$. Combining this observation with (2.10) we see that the quantity $x(n) * g^*(-n)$ can be computed using the schematic shown in Fig. 3. In the figure, the "correlator" computes $X_k(z^M) \tilde{S}_k(z)$.

Notice that unlike traditional convolution theorems, we do not have to apply an "inverse transform" after performing the transform domain operations. This is true even if the filter bank is the DFT filter bank (i.e., Fig. 2 with $E(z)$ equal to the DFT matrix). So even in the special case of DFT filter banks, the above result is fundamentally different from well-known DFT based convolutions. See Section VII-A for further remarks about this.

Comments on complexity: Computational complexity is not the main advantage of the method of subband convolution. Assume for simplicity that $x(n)$ and $g(n)$ are N -point sequences. Then direct convolution of $x(n)$ and $g^*(-n)$ (without using standard fast techniques) requires N^2 multipliers. Assuming that N is much larger than the lengths of the subband filters $H_k(z)$ (so that the multiplications required to implement analysis filters are negligible) the signals $x_k(n)$ and $g_k(n)$ have lengths $\approx N/M$. Each subband convolution requires nearly $(N/M)^2$ multiplications, so that the total number of multiplications for all the M values of i in (2.9) is nearly N^2 again. It is true that we can employ the FFT, or even the fast short convolution algorithms in the subbands, but again this is not the main point of the discussion.

The above reasoning does not hold if the analysis filters have length comparable to those of $x(n)$ and $g(n)$. In this case, the complexity of the analysis bank becomes comparable to the direct convolution of $x(n)$ with $g(n)$ and this additional overhead may overshadow the advantages of subband convolution: Recall that the actual advantage of the (paraunitary) subband convolver is that it allows us to allocate the computational accuracy (i.e., bits) among the subbands, resulting in a coding gain as elaborated in Section III. In fact, considerable coding gain can be obtained even in the special case where the analysis filters have small length (e.g., $\leq M$), as discussed in Sections IV and V.

B. Orthonormal Filter Bank with Unequal Decimation Ratios

Now consider the case where the decimation ratios n_k in Fig. 1 are possibly unequal integers such that

$$\sum_{k=0}^{M-1} \frac{1}{n_k} = 1. \quad (2.11)$$

This condition implies that we have a maximally decimated system. The design of such nonuniform systems has received attention recently [21]–[23]. Such a system can be regarded as a discrete time wavelet decomposition system. The analysis bank is the "wavelet transformer"

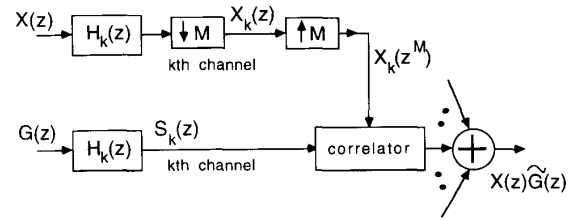


Fig. 3. Implementation of the paraunitary convolver.

and the synthesis bank the inverse transformer. Assuming perfect reconstruction (i.e., $\hat{x}(n) = x(n)$) we can express the signal $x(n)$ in terms of the synthesis filters $F_k(z)$ and the wavelet coefficients $X_k(z)$ as follows:

$$X(z) = \sum_{k=0}^{M-1} F_k(z) X_k(z^{n_k}) \quad (2.12a)$$

i.e., in the time domain,

$$x(n) = \sum_{k=0}^{M-1} \sum_l x_k(l) f_k(n - n_k l). \quad (2.12b)$$

The doubly indexed set of sequences

$$\xi_{k,l}(n) \triangleq f_k(n - n_k l), \quad 0 \leq k \leq M-1, \\ -\infty \leq l \leq \infty \quad (2.13)$$

are therefore the basis functions for the expansion of $x(n)$. Note that the sequence $x(n)$ and the basis functions $\xi_{k,l}(n)$ are, in general, doubly infinite in extent (i.e., $-\infty \leq n \leq \infty$). Special cases of this system based on binary and nonbinary tree structures (wavelet packets) have been reported earlier [5], [7], [10], [24], [25].

1) *Orthonormality (Nonuniform Case):* The above basis is said to be orthonormal if

$$\sum_n \underbrace{f_k(n - n_k l)}_{\xi_{k,l}(n)} \underbrace{f_m^*(n - n_m i)}_{\xi_{m,i}^*(n)} = \delta(k - m) \delta(l - i). \quad (2.14)$$

In terms of the synthesis filters, the orthonormality property is

$$\sum_n f_k(n) f_m^*(n + n_k l - n_m i) = \delta(k - m) \delta(l - i). \quad (2.15)$$

This is a generalization of the orthonormality property (2.5) which followed earlier from paraunitariness. Let $n_{k,m}$ denote the greatest common divisor of n_k and n_m , i.e.,

$$n_{k,m} = \gcd(n_k, n_m). \quad (2.16)$$

We can then rewrite (2.15) as [24]

$$\sum_n f_k(n) f_m^*(n + n_{k,m} p) = \delta(k - m) \delta(p) \quad (2.17)$$

(see Appendix A). In the z -domain this is equivalent to

$$(F_k(z) \tilde{F}_m(z))_{n_{k,m}} = \delta(k - m) \quad (\text{orthonormality}). \quad (2.18)$$

A simple example of a perfect reconstruction orthonormal filter bank with unequal n_k is obtained by use of a binary

tree structure [5] with paraunitary polyphase matrices at each level [10], [12], [24]. This results in filter responses that have an octave spacing. For a preview, the reader can see Fig. 7. The heights of the filters are unequal because the energy of each filter has to be unity (as seen by setting $k = m$ and $p = 0$ in (2.17)).

Properties of nonuniform orthonormal filter banks: Some crucial features of maximally decimated orthonormal filter banks are summarized next. These are elaborated in Appendix B.

1) For a perfect reconstruction system (i.e., $\hat{x}(n) = x(n)$ in Fig. 1(a)), the analysis and synthesis filters are related as $F_k(z) = \tilde{H}_k(z)$. This implies, in particular, that $|H_k(e^{j\omega})| = |F_k(e^{j\omega})|$. As noted above, all the filters have unit energy, that is

$$\int_0^{2\pi} |H_k(e^{j\omega})|^2 d\omega/2\pi = \int_0^{2\pi} |F_k(e^{j\omega})|^2 d\omega/2\pi = 1.$$

2) Since $F_k(z) = \tilde{H}_k(z)$, the analysis filters also satisfy orthonormality, i.e.,

$$(H_k(z)\tilde{H}_m(z))_{|n_{k,m}} = \delta(k - m).$$

3) The filters satisfy a generalized version of the power complementary property, viz.,

$$\sum_{k=0}^{M-1} \frac{|H_k(e^{j\omega})|^2}{n_k} = 1, \text{ and } \sum_{k=0}^{M-1} \frac{|F_k(e^{j\omega})|^2}{n_k} = 1. \quad (2.19)$$

4) Let L denote the least common multiple (lcm) of the decimation ratios. Then orthonormality also implies $(F_k(z)\tilde{F}_m(z))_{|L} = \delta(k - m)$. This means that the $L \times M$ matrix $F(z)$ with elements

$$[F(z)]_{n,i} = F_i(zW_L^n), \quad 0 \leq i \leq M-1, \\ 0 \leq n \leq L-1$$

is paraunitary, that is $\tilde{F}(z)F(z) = LI$. So the orthonormality of a nonuniform filter bank is essentially the paraunitary property in disguise. In fact, it has been observed [21], [22] that the nonuniform system can be redrawn as a uniform L -band filter bank (with L -fold decimators); this "bigger system" is paraunitary if and only if the smaller nonuniform system is orthonormal (Appendix B).

2) **Derivation of the Convolution Theorem (Nonuniform Case):** Assume that we have perfect reconstruction, i.e., $\hat{X}(z) = X(z)$ and $\hat{G}(z) = G(z)$. Using the expression (2.12a) for $X(z)$ and similarly for $G(z)$, we have

$$X(z)\tilde{G}(z) = \sum_{k=0}^{M-1} \sum_{m=0}^{M-1} F_k(z)\tilde{F}_m(z)X_k(z^{n_k})\tilde{G}_m(z^{n_m}). \quad (2.20)$$

Let L be the least common multiple of the decimation ratios, i.e.,

$$L = \text{lcm}\{n_k\}. \quad (2.21)$$

For $0 \leq k, m \leq M-1$ we then have

$$L = n_k p_k, \quad L = n_m p_m \quad (2.22)$$

for some integers p_k and p_m . Consider now the L -fold decimated version of $X(z)\tilde{G}(z)$. Using the above decomposition of L and the identity (1.1), we can write

$$(X(z)\tilde{G}(z))_{|L} = \sum_{k=0}^{M-1} \sum_{m=0}^{M-1} ((F_k(z)\tilde{F}_m(z))_{|n_{k,m}} X_k(z^{n_k/n_{k,m}}) \cdot \tilde{G}_m(z^{n_m/n_{k,m}}))_{|p_{k,m}} \quad (2.23)$$

since $n_{k,m}$ is a common factor of n_k and n_m . Using the orthonormality property (2.18) this simplifies to

$$(X(z)\tilde{G}(z))_{|L} = \sum_{k=0}^{M-1} (X_k(z)\tilde{G}_k(z))_{|p_k}. \quad (2.24a)$$

Equivalently, in the time domain

$$(x(n) * g^*(-n))_{|L} = \sum_{k=0}^{M-1} (x_k(n) * g_k^*(-n))_{|p_k}. \quad (2.24b)$$

Again, there is no cross-coupling between subbands. To obtain all the samples of the convolution $x(n) * g^*(-n)$, we have to (at least conceptually) repeat the above with the shifted versions $g(n-i)$, $0 \leq i \leq L-1$, even though a simpler procedure will be described below. The main result is summarized as follows.

Theorem 2.2: Convolution theorem for orthonormal nonuniform filter banks. Consider the maximally decimated filter bank of Fig. 1, and ignore the quantizers Q_k . Let

$$L = \text{lcm}\{n_i\}, \quad n_{k,m} = \text{gcd}(n_k, n_m), \quad \text{and} \quad p_k = L/n_k. \quad (2.25)$$

Assume that the system has perfect reconstruction ($\hat{x}(n) = x(n)$ for any $x(n)$) and that the synthesis filters are orthonormal, i.e., satisfy (2.17) or equivalently (2.18). Then the L -fold decimated version of the convolution $x(n) * g^*(-n)$ can be computed by computing the p_k -fold decimated versions of the convolutions $x_k(n) * g_k^*(-n)$, and adding them. We can obtain all the samples of the convolution by repeating this for L successively shifted versions of $g(n)$. \diamond

Comments

1) **Implementation that obtains all samples.** Let $s(n)$ denote the (undecimated) output of $H_k(z)$ in response to the unshifted input $g(n)$. Then

$$S(z) = \sum_{m=0}^{n_k-1} z^m G_k^{(m)}(z^{n_k}).$$

Thus, we can recover all $G_k^{(m)}(z)$ from the undecimated signal $S(z)$. We can obtain an efficient implementation of the convolution as follows: from Appendix B-3 we know that the given filter bank is equivalent to an L -channel uniform filter bank with equal decimation ratio L in all channels. For this uniform system, Theorem 2.1 holds. We can therefore obtain an implementation of $X(z)\tilde{G}(z)$ by using the scheme of Fig. 3, with M replaced by L and the filters $\{H_k(z)\}$ replaced by $\{H'_k(z)\}$ as described in Appendix B-3.

2) *Complexity*. From a computational complexity viewpoint, the comments following (2.10) continue to hold. It can be shown that the number of multiplications for a direct convolution $x(n) * g^*(-n)$ are nearly the same as the total number of multiplications required to perform all the necessary subband convolutions. (This neglects the multiplications required to implement the analysis filters $H_k(z)$ and assumes that the lengths of $H_k(z)$ are much smaller than those of $x(n)$ and $g(n)$.)

3) *Parseval's relation*. If we evaluate the 0th sample of the convolution (2.24b) we obtain

$$\sum_n x(n)g^*(n) = \sum_{k=0}^{M-1} \sum_n x_k(n)g_k^*(n). \quad (2.26)$$

This can be regarded as the equivalent of Parseval's relation [13], in the world of (nonuniform) orthonormal filter bank transforms. With $x(n) = g(n)$, this reduces to the energy balance equation

$$\sum_n |g(n)|^2 = \sum_{k=0}^{M-1} \sum_n |g_k(n)|^2 \quad (\text{Energy conservation}). \quad (2.27)$$

III. CODING GAIN OF ORTHONORMAL CONVOLVERS

Fig. 1 shows the paraunitary convolver with quantizers inserted in the subbands of $x(n)$. We will first consider the uniform case ($n_k = M$ for all k). The nonuniform case will be addressed in Section III-C. Assume that $g(n)$ is a fixed filter with no quantizers in its subbands. (This assumption can be removed, but only with considerable loss of simplicity of mathematics.) For simplicity of analysis we assume that $x(n)$, $g(n)$ and the filter coefficients in $H_k(z)$ are real so that $x_k(n)$ and $g_k(n)$ are real. This enables us to deal with quantizers that operate on real inputs.

Let b_k denote the number of bits per sample of $x_k(n)$, permitted by the quantizer Q_k . Thus the average bit rate is

$$b = \frac{1}{M} \sum_{k=0}^{M-1} b_k \quad (3.1)$$

i.e., on the average, we have used b bits per sample of $x(n)$.

Because of the quantization in the subbands, the output of the paraunitary convolver is different from the ideal result $x(n) * g^*(-n)$. To analyze this error, we replace the quantizers Q_k with the noise sources $q_k(n)$ as shown in Fig. 4. Consider the paraunitary convolution formula (2.8). In the presence of quantizers, we are actually computing

$$\sum_{k=0}^{M-1} (x_k(n) + q_k(n)) * g_k^*(-n). \quad (3.2)$$

(According to the realness assumption the conjugate sign is redundant, but we show it for consistency with previous sections.) The quantization error is therefore

$$q(n) = \sum_{k=0}^{M-1} q_k(n) * g_k^*(-n). \quad (3.3)$$

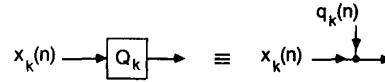


Fig. 4. The quantizer and its noise model.

The noise model: To perform a statistical analysis, we will make the following assumptions:

1) $x(n)$ is a zero-mean wide-sense stationary random process so that the subband signals $x_k(n)$ are zero-mean WSS with some variance, say, $\sigma_{x_k}^2$. We consider $g(n)$ to be deterministic sequence.

2) The quantizer noise source $q_k(n)$ is zero-mean and white, with variance $\sigma_{q_k}^2$. Also $q_k(n)$ is uncorrelated to $q_m(i)$, $k \neq m$, and to the input $x(n)$ (hence to the quantizer input $x_k(n)$).

It should be noticed that the above assumptions are reasonable as long as the bit rates b_k are moderate or high [26]. In any case, in the absence of such assumptions, it is not usually possible to find an expression for error variance. (However, in the special case of subband coders (where $g(n) = \delta(n)$), these assumptions can be relaxed. See Section III-D and Appendix C.)

A. Expression for the Error Variance

Let $\sigma_{x_k}^2$ denote the variance of $x_k(n)$, and $\sigma_{q_k}^2$ the variance of the quantizer error $q_k(n)$. In order to equalize the overflow probability across all the M subbands, these two should be related as

$$\sigma_{q_k}^2 = c \sigma_{x_k}^2 2^{-2b_k}. \quad (3.4)$$

(For a detailed explanation of this equation, see, e.g., [17, ch. 4], or [12, appendix, sec. C.1].) Here c is a constant, identical for all subbands (which is a valid assumption if all $x_k(n)$ have similar type of distribution, e.g., all Gaussian).

Using (3.3) and the noise model assumptions stated earlier, the variance of $q(n)$ can be expressed as

$$\begin{aligned} \sigma_{q(n)}^2 &= \sum_{k=0}^{M-1} \sigma_{q_k}^2 \sum_l |g_k(l)|^2 \\ &= c \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \sum_l |g_k(l)|^2 \quad (\text{from (3.4)}). \end{aligned} \quad (3.5)$$

This is for $i = 0$ in (2.9). For arbitrary i , the filter $g(n)$ is replaced with $g(n - i)$, and the above equation is modified to

$$\begin{aligned} \sigma_{q(n-i)}^2 &= c \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \sum_l |g_k^{(i)}(l)|^2, \\ 0 \leq i \leq M-1 \end{aligned} \quad (3.6)$$

where $g_k^{(i)}(n)$ is the k th subband output in response to $g(n - i)$. The dependence on i is removed by averaging over all i . The resulting average variance of $q(n)$ is given by

$$\sigma_{q,PU}^2 = \frac{c}{M} \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \alpha_k^2 \quad (3.7)$$

where

$$\alpha_k^2 \triangleq \sum_{i=0}^{M-1} \sum_l |g_k^{(i)}(l)|^2. \quad (3.8)$$

The inner summation above represents the energy in the k th subband in response to the input $g(n-i)$. The outer summation removes the dependence on i . Thus α_k^2 is proportional to the average energy of $g(n)$ in the k th subband. Using the paraunitary property, it can be shown that $\sum_k \alpha_k^2/M$ is the total energy in $g(n)$ (see (3.22) later).

The PU in the subscript in (3.7) is a reminder of paraunitary. Equation (3.7) gives the average error variance (over a period of length M), and is independent of time.

B. Coding Gain of the Paraunitary Convolver

Now consider direct convolution $x(n) * g^*(-n)$. Suppose $x(n)$ is directly quantized to b bits before convolution. Denoting $e(n)$ as the quantization error, the result of quantization is $[x(n) + e(n)] * g^*(-n)$ so that the error is $e(n) * g^*(-n)$. Under usual noise model assumptions, the variance of this error is

$$\sigma_{q,\text{direct}}^2 = \sigma_e^2 \sum_n |g(n)|^2 \quad (3.9)$$

where σ_e^2 is the variance of $e(n)$, which can be expressed, similar to (3.4), as $\sigma_e^2 = c\sigma_x^2 2^{-2b}$, where σ_x^2 is the variance of $x(n)$. Thus

$$\sigma_{q,\text{direct}}^2 = c\sigma_x^2 2^{-2b} \sum_n |g(n)|^2. \quad (3.10)$$

The ratio

$$G_{\text{PU}}(M) = \frac{\sigma_{q,\text{direct}}^2}{\sigma_{q,\text{PU}}^2} \quad (3.11)$$

is the coding gain of the paraunitary convolver. The argument M is a reminder that there are M subbands in the system. Substituting from (3.7) and (3.10), this becomes

$$G_{\text{PU}}(M) = \frac{2^{-2b} \sigma_x^2 \sum_n |g(n)|^2}{\frac{1}{M} \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \alpha_k^2}. \quad (3.12)$$

In this expression, $\sigma_{x_k}^2$ is the variance of the k th subband signal derived from the input $x(n)$, and $\alpha_k^2 \geq 0$ is related to the k th subband of the filter $g(n)$. And b is the average bit rate (3.1). Notice that $\sigma_{x_k}^2$ and α_k^2 depend on the analysis filter response $H_k(e^{j\omega})$.

1) *Optimum Bit Allocation:* Under the average bit-rate constraint (3.1), we can maximize the coding gain by optimally allocating the bits b_k among subbands. The idea is very similar to the counterpart in subband coding [17]. For this we note that the numerator of (3.12) is independent of the bit allocation. We only have to minimize the denominator. For this we invoke the arithmetic-geometric

mean inequality [27] (AM-GM inequality) which states this: if P_k , $0 \leq k \leq M-1$, is a set of nonnegative numbers, then

$$\frac{1}{M} \sum_{k=0}^{M-1} P_k \geq \left(\prod_{k=0}^{M-1} P_k \right)^{1/M} \quad (3.13)$$

with equality if and only if $P_k = P$ for all k . Using this in conjunction with (3.1) we can show that

$$\frac{1}{M} \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \alpha_k^2 \geq 2^{-2b} \prod_{k=0}^{M-1} (\sigma_{x_k}^2 \alpha_k^2)^{1/M} \quad (3.14)$$

with equality if and only if all terms on the left side above are equal. Since the quantizer variances $\sigma_{q_k}^2$ are given by (3.4), we see that the above condition for equality implies

$$\sigma_{q_k}^2 = c\sigma_{x_k}^2 2^{-2b_k} = \frac{\text{constant}}{\alpha_k^2}.$$

The output noise variance due to the k th quantizer (k th term in (3.7)) is therefore independent of k .

We obtain the formula for the optimal bit allocation by setting all the terms on the left side of (3.14) to be equal. The result is

$$b_k = C + \frac{1}{2} \log_2 (\sigma_{x_k}^2 \alpha_k^2) \quad (3.15)$$

for some C . By using (3.1) we can eliminate C and obtain

$$b_k = b + 0.5 \log_2 (\sigma_{x_k}^2 \alpha_k^2) - \frac{0.5}{M} \sum_{i=0}^{M-1} \log_2 (\sigma_{x_i}^2 \alpha_i^2). \quad (3.16)$$

This is very similar to the expressions which can be found in [17], [18] for traditional subband coding systems. The difference is that the product $\sigma_{x_k}^2 \alpha_k^2$ appears in the place of $\sigma_{x_k}^2$. Thus, the energy of the signal as well as the filter $g(n)$ in the k th subband determine the bits b_k . (These formulas are similar to the case of subband coding with frequency weighting; see [17, p. 532].) For high bit rate coding, the above expression is useful. As in subband coding, b_k might turn out to be nonintegral, and sometimes negative if b is not large enough.

2) *Optimum Coding Gain:* The optimum convolutional coding gain is obtained when equality holds in (3.14), i.e., when all the terms on the left side of (3.14) are equal. The optimum value is

$$G_{\text{PU,optimal}}(M) = \frac{\sigma_x^2}{\left(\prod_{k=0}^{M-1} \sigma_{x_k}^2 \right)^{1/M}} \times \frac{\sum_n |g(n)|^2}{\left(\prod_{k=0}^{M-1} \alpha_k^2 \right)^{1/M}}. \quad (3.17)$$

Notice that the above analysis holds for any filter-bank convolver with paraunitary polyphase matrix, regardless of the quality of the filter responses. The filter responses will in turn determine the values of $\sigma_{x_k}^2$ and α_k^2 for fixed $g(n)$ and $x(n)$.

Lemma 3.1: $G_{\text{PU,optimal}}(M) \geq 1$ regardless of the choice of paraunitary filters $H_k(z)$. Moreover,

$G_{\text{PU, optimal}}(M) = 1$ if and only if the subband variance $\sigma_{x_k}^2$ and the quantity α_k^2 are independent of k . \diamond

Proof: We will rewrite the optimal coding gain (3.17) by expressing σ_x^2 in terms of $\sigma_{x_k}^2$ and $\sum_n |g(n)|^2$ in terms of α_k^2 .

The variance of the output of $H_k(z)$ is also the variance of the decimated subband signal $x_k(n)$ so that

$$\sigma_{x_k}^2 = \frac{1}{2\pi} \int_0^{2\pi} S_{xx}(e^{j\omega}) |H_k(e^{j\omega})|^2 d\omega \quad (3.18)$$

where $S_{xx}(e^{j\omega})$ is the power spectral density of $x(n)$. The paraunitary property $\tilde{E}(z)E(z) = I$ implies

$$\sum_{k=0}^{M-1} |H_k(e^{j\omega})|^2 = M.$$

By computing $\sum_k \sigma_{x_k}^2$ from (3.18) we therefore obtain

$$\frac{1}{M} \sum_{k=0}^{M-1} \sigma_{x_k}^2 = \sigma_x^2. \quad (3.19)$$

Next consider the signals generated by the filter bank in response to $g(n-i)$ (Fig. 5). Define the vectors

$$\hat{\mathbf{g}}^{(i)}(n) = \begin{bmatrix} g(Mn-i) \\ g(Mn-1-i) \\ \vdots \\ g(Mn-M+1-i) \end{bmatrix},$$

$$\mathbf{g}^{(i)}(n) = \begin{bmatrix} g_0^{(i)}(n) \\ g_1^{(i)}(n) \\ \vdots \\ g_{M-1}^{(i)}(n) \end{bmatrix} \quad (3.20)$$

for $0 \leq i \leq M-1$. The superscript i is a reminder that the input is $g(n-i)$. Using the paraunitary property, we conclude [12]

$$\sum_n [\hat{\mathbf{g}}^{(i)}(n)]^\dagger \hat{\mathbf{g}}^{(i)}(n) = \sum_n [\mathbf{g}^{(i)}(n)]^\dagger \mathbf{g}^{(i)}(n). \quad (3.21)$$

The left-hand side is the energy $\sum_n |g(n)|^2$. Combining this with the definition (3.8) of α_k^2 , we obtain

$$\sum_n |g(n)|^2 = \frac{1}{M} \sum_{k=0}^{M-1} \alpha_k^2. \quad (3.22)$$

Substituting (3.19) and (3.22) into (3.17), we arrive at

$$G_{\text{PU, optimal}}(M) = \frac{\frac{1}{M} \sum_{k=0}^{M-1} \sigma_{x_k}^2}{\left(\prod_{k=0}^{M-1} \sigma_{x_k}^2 \right)^{1/M}} \times \frac{\frac{1}{M} \sum_{k=0}^{M-1} \alpha_k^2}{\left(\prod_{k=0}^{M-1} \alpha_k^2 \right)^{1/M}}. \quad (3.23)$$

Using the arithmetic-geometric mean inequality we conclude that $G_{\text{PU, optimal}}(M) \geq 1$, with equality if and only if $\sigma_{x_k}^2$ and α_k^2 are independent of k . $\nabla\nabla\nabla$

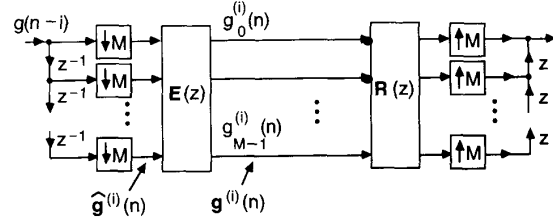


Fig. 5. Response of the filter bank to a shifted input.

Thus, in general, the gain can exceed unity for two possible reasons. First the subband variances $\sigma_{x_k}^2$ could be different for different k . And second, the quantity α_k may not be the same in all subbands.

Notice that the above proof uses the paraunitary property. The property $G_{\text{PU, optimal}}(M) \geq 1$ cannot be claimed for a convolver based on an arbitrary filter bank (i.e., without paraunitary property). The appearance of the arithmetic-geometric mean ratio in the coding gain has been observed in other contexts in traditional subband coding applications. It has been formally proved for the case of ideal brick-wall filters and for the case of orthogonal transform coding [17]. Such an expression has also been used for other types of (nonideal) filter banks [28]. The true justification for such use is based on the paraunitary property, as shown above and in [29].

Special Cases: Paraunitary filter banks are special cases of perfect reconstruction filter banks [2], [3]. However, they cover a wide range of practical filter banks. In fact, some of the approximate reconstruction systems (viz., the pseudo-QMF banks [30]–[33]) are known to satisfy the paraunitary property “approximately” (see [34]), even though these approximate systems were developed before paraunitary filter banks were reported.

1) A special case of paraunitary systems, primarily of theoretical interest, arises when the filters $H_k(e^{j\omega})$ are equispaced ideal brick-wall filters as shown in Fig. 6. In this case

$$F_k(e^{j\omega}) = H_k(e^{j\omega}) = \begin{cases} \sqrt{M} & \text{if } \omega \in k\text{th passband} \\ 0 & \text{otherwise} \end{cases} \quad (3.24)$$

and it can be shown that $E(e^{j\omega})$ is paraunitary (see [12, sec. 6.2.2]). In this case, we have

$$\sigma_{x_k}^2 = M \int_{k\text{th band}} S_{xx}(e^{j\omega}) d\omega / 2\pi \quad (3.25)$$

where $S_{xx}(e^{j\omega})$ is the power spectrum of $x(n)$. The system (3.24) will be called the ideal subband convolver (SBC).

2) A second special case of theoretical interest arises when $H_k(z) = z^{-k}$ for all k . In this case the above results still hold (since $E(z) = I$ which is paraunitary); and the coding gain can be verified to be unity.

3) *Case of white input.* Suppose $x(n)$ is zero mean and white. Then $\sigma_{x_k}^2$ is identical for all k . This follows because paraunitariness implies in particular, that the energy

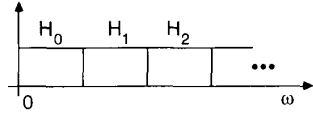


Fig. 6. Magnitude response of ideal brick-wall analysis filters. Synthesis filters for perfect reconstruction have the same magnitude responses.

$\int_0^{2\pi} |H_k(e^{j\omega})|^2 d\omega/2\pi$ is identical for all k (Appendix B). In this case, the coding gain can still exceed unity, because α_k^2 may not be identical for all k .

C. Coding Gain for the Nonuniform Orthonormal Convolver

In the nonuniform case, (2.24) gives the L -fold decimated version of the convolution. To obtain all samples of the convolution, we repeat this operation with $g(n)$ replaced by $g(n-i)$, i.e., $g_k(n)$ replaced by $g_k^{(i)}(n)$ for $0 \leq i \leq L-1$.

With quantizers inserted as in Fig. 1(a), we can replace them with noise sources $q_k(n)$ as in Fig. 4. With $x(n)$ and $g(n-i)$ used as the filter bank inputs, the error in the computation of $\sum_{k=0}^{M-1} x_k(n) * [g_k^{(i)}(-n)]^*$ is therefore $\sum_{k=0}^{M-1} q_k(n) * [g_k^{(i)}(-n)]^*$. Proceeding as before, we find the variance of this error to be

$$\sum_{k=0}^{M-1} \sigma_{q_k}^2 \sum_n |g_k^{(i)}(n)|^2. \quad (3.26)$$

Averaging over the L values of i , we obtain the average variance of the error $q(n)$ in the convolution as

$$\sigma_{q,\perp}^2 = \frac{1}{M} \sum_{k=0}^{M-1} \sigma_{q_k}^2 \alpha_k^2 = \frac{c}{M} \sum_{k=0}^{M-1} 2^{-2b_k} \sigma_{x_k}^2 \alpha_k^2, \quad (3.27)$$

where

$$\alpha_k^2 \triangleq (M/L) \sum_{i=0}^{L-1} \sum_n |g_k^{(i)}(n)|^2. \quad (3.28)$$

Here we have used (3.4). The quantity $\sigma_{q,\perp}^2$ is the “output error variance” of the convolver. The subscript \perp stands for “orthonormal” filter banks.

The bit rate for the k th subband is b_k/n_k . Assume that the total bit rate is constrained to be b . Then the bit rate constraint is

$$\sum_{k=0}^{M-1} \frac{b_k}{n_k} = b. \quad (3.29)$$

To obtain the optimal bit allocation, we can minimize $\sigma_{q,\perp}^2$ under the above constraint by use of the Lagrange multiplier method. That is, form the Lagrangian $\phi = \sigma_{q,\perp}^2 - \lambda(\sum_{k=0}^{M-1} b_k/n_k - b)$ and set $\partial\phi/\partial b_k = 0$. This results in the set of equations

$$2^{2b_k} = D \sigma_{x_k}^2 \alpha_k^2 n_k, \quad 0 \leq k \leq M-1 \quad (3.30)$$

where D is a constant independent of k .³ Taking logarithm

³The fact that this represents a minimum rather than maximum can be verified in many ways. For example, one can verify in this case that the Hessian of the Lagrangian [35] is a diagonal matrix with positive elements.

and using (3.29), we can evaluate the constant D to be

$$D = \frac{2^{2b}}{\prod_{i=0}^{M-1} (\sigma_{x_i}^2 \alpha_i^2 n_i)^{1/n_i}}. \quad (3.31)$$

Substituting this into (3.30) and taking logarithm, we obtain the following formula:

$$b_k = b + 0.5 \log_2 (n_k \sigma_{x_k}^2 \alpha_k^2) - 0.5 \sum_{i=0}^{M-1} \frac{\log_2 (n_i \sigma_{x_i}^2 \alpha_i^2)}{n_i} \quad (3.32)$$

for optimal bit allocation. Under this condition, the variance of the k th quantizer noise is given by

$$\sigma_{q_k}^2 = c 2^{-2b_k} \sigma_{x_k}^2 = \frac{c}{D \alpha_k^2 n_k} = \frac{c 2^{-2b} \prod_{i=0}^{M-1} (\sigma_{x_i}^2 \alpha_i^2 n_i)^{1/n_i}}{\alpha_k^2 n_k} \quad (3.33)$$

which is proportional to $1/(\alpha_k^2 n_k)$. With optimum bit allocation, the output noise variance contributed by the k th quantizer (k th term in (3.27)) simplifies to $c/(DMn_k)$, and is proportional to $1/n_k$. The total output noise variance is

$$\sigma_{q,\perp}^2 = \frac{c}{DM} \sum_{k=0}^{M-1} 1/n_k = \frac{c}{DM} = \frac{c 2^{-2b}}{M} \prod_{i=0}^{M-1} (\sigma_{x_i}^2 \alpha_i^2 n_i)^{1/n_i}. \quad (3.34)$$

The convolutional coding gain, defined as $G_{\perp, \text{optimal}}(M) = \sigma_{q, \text{direct}}^2 / \sigma_{q,\perp}^2$ can now be calculated. Thus using (3.10) and (3.34) we obtain

$$G_{\perp, \text{optimal}}(M) = \frac{\sigma_x^2 \sum_n |g(n)|^2}{\frac{1}{M} \prod_{i=0}^{M-1} (n_i \sigma_{x_i}^2 \alpha_i^2)^{1/n_i}}. \quad (3.35)$$

Notice that these results reduce to those in Section III-B if we set $n_k = M$ for all k . Another special case of interest in many applications (speech and image coding) is the filter bank with analysis filter responses resembling the one in Fig. 7. The responses have an octave spacing and correspondingly increasing bandwidths (constant Q filter bank). (The unequal filter heights are such that all filters have the same energy.) Such a system can be generated by use of a tree-structured system, where one of the two signals from the previous stage is further split into two in the next stage [5], and so forth. The orthonormality property can be satisfied in such a system by use of 2×2 paraunitary polyphase matrices at each level of the tree. The above theory can be applied for these systems, with

$$n_0 = n_1 = 2n_2 = 4n_3 = \dots$$

Lemma 3.2: $G_{\perp, \text{optimal}}(M) \geq 1$ regardless of the choice of the orthonormal filters $H_k(z)$. Moreover $G_{\perp, \text{optimal}}(M) = 1$ if and only if $\sigma_{x_k}^2$ and $n_k \alpha_k^2$ are independent of k . \diamond

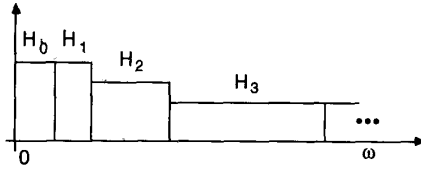


Fig. 7. Magnitude responses of ideal analysis filters, for a well-known class of nonuniform filter banks.

Proof⁴: We will first express the numerator in (3.35) in terms of the quantities in the denominator. By using (3.18), and the property (2.19), we can deduce that

$$\sigma_x^2 = \sum_{k=0}^{M-1} \frac{\sigma_{xk}^2}{n_k}. \quad (3.36)$$

Furthermore, the total energy of $g(n)$ can be expressed as

$$\sum_n |g(n)|^2 = \frac{1}{M} \sum_{k=0}^{M-1} \alpha_k^2. \quad (3.37)$$

(For this, just apply Parseval's theorem (2.27). With $g(n)$ replaced by $g(n-i)$ the left side of (2.27) is unchanged but $g_k(n)$ gets replaced with $g_k^{(i)}(n)$. Using the definition (3.28) of α_k^2 , we obtain the desired result.) Using the above two expressions, we obtain

$$G_{\perp, \text{optimal}}(M) = \left(\frac{\sum_{k=0}^{M-1} \sigma_{xk}^2 / n_k}{\prod_{k=0}^{M-1} (\sigma_{xk}^2)^{1/n_k}} \right) \cdot \left(\frac{\sum_{k=0}^{M-1} \alpha_k^2}{\prod_{k=0}^{M-1} (n_k \alpha_k^2)^{1/n_k}} \right). \quad (3.38)$$

To prove that this is ≥ 1 , we rewrite

$$\sum_{k=0}^{M-1} \frac{\sigma_{xk}^2}{n_k} = \sum_{k=0}^{M-1} \frac{p_k \sigma_{xk}^2}{L} \quad (3.39)$$

where L is the lcm of $\{n_k\}$ and $L = p_k n_k$, as before. Since $\sum_k p_k = \sum_k (L/n_k) = L$, the right-hand side in (3.39) is really a sum of L nonnegative terms (σ_{x0}^2 occurs p_0 times, σ_{x1}^2 occurs p_1 times and so forth). Applying the arithmetic-geometric mean inequality and simplifying, we conclude that

$$\sum_{k=0}^{M-1} \frac{\sigma_{xk}^2}{n_k} \geq \prod_{k=0}^{M-1} (\sigma_{xk}^2)^{1/n_k} \quad (3.40)$$

with equality if and only if all the σ_{xk}^2 are identical. Next, we can write

$$\sum_{k=0}^{M-1} \alpha_k^2 = \frac{1}{L} \sum_{k=0}^{M-1} p_k (n_k \alpha_k^2) \quad (\text{since } p_k n_k = L). \quad (3.41)$$

⁴This proof was suggested by T. Chen, graduate student, California Institute of Technology.

As before, the right-hand side is a sum of L nonnegative terms, and we obtain after some simplification

$$\sum_{k=0}^{M-1} \alpha_k^2 \geq \prod_{k=0}^{M-1} (n_k \alpha_k^2)^{1/n_k} \quad (3.42)$$

with equality if and only if $n_k \alpha_k^2$ has the same value for all k . Using these inequalities in (3.38), the lemma follows immediately. $\nabla \nabla \nabla$

D. The Special Case of Traditional Subband Coding ($g(n) = \delta(n)$)

The results derived above for the paraunitary convolvers (uniform as well as nonuniform) can be used to derive the optimal bit allocation and coding gain for orthonormal subband coding systems, i.e., systems of the form in Fig. 1(a). This is done by setting $g(n) = \delta(n)$. Under this condition, the quantity $g_k^{(i)}(n)$ is the decimated impulse response $h_k(nn_k - i)$, where $h_k(n)$ is the impulse response of the analysis filter $H_k(z)$. Using the fact that the analysis filters have unit energy under the orthonormality constraint, one can verify that $\alpha_k^2 = M/n_k$. Substituting this we obtain the reconstruction error variance, i.e., variance of $x(n) - \hat{x}(n)$ in Fig. 1(a). This can be obtained from (3.27) as

$$\sigma_{q, \perp}^2 = \sum_{k=0}^{M-1} \frac{\sigma_{qk}^2}{n_k}. \quad (3.43)$$

The optimal bit allocation rule (which minimizes the above expression) reduces to

$$b_k = b + 0.5 \log_2 (\sigma_{xk}^2) - 0.5 \sum_{i=0}^{M-1} \frac{\log_2 (\sigma_{xi}^2)}{n_i} \quad (3.44)$$

and the optimized coding gain becomes

$$G_{\perp, \text{optimal}}(M) = \frac{\sigma_x^2}{\prod_{k=0}^{M-1} (\sigma_{xk}^2)^{1/n_k}} \quad (3.45)$$

and can be rewritten, using the orthonormality of the analysis bank (as in the proof of Lemma 3.2), as

$$G_{\perp, \text{optimal}}(M) = \frac{\sum_{k=0}^{M-1} \sigma_{xk}^2 / n_k}{\prod_{k=0}^{M-1} (\sigma_{xk}^2)^{1/n_k}}. \quad (3.46)$$

Clearly, $G_{\perp, \text{optimal}}(M) \geq 1$ with equality if and only if σ_{xk}^2 is the same for all k (from Lemma 3.2). Since $\alpha_k^2 = M/n_k$ in this case, we see from (3.33) that the variance σ_{qk}^2 of the k th quantizer noise source is independent of k under optimal bit allocation, and is given by

$$\sigma_{qk}^2 = c 2^{-2b} \prod_{i=0}^{M-1} (\sigma_{xi}^2)^{1/n_i}. \quad (3.47)$$

However, the contribution to the output noise variance $\sigma_{q, \perp}^2$, coming from the k th quantizer (k th term in (3.43), is proportional to $1/n_k$.

Summarizing, the above expressions are applicable to any subband coder (possibly unequal decimation ratios, but maximally decimated) with orthonormal filters, under the noise model assumptions stated at the beginning of Section III. Also see [29] for more details. The further special case where $n_k = M$ has been reported in many references in the past [17], [18], [28].

Some subtleties about basic assumptions: The above references assume ideal nonoverlapping subband filters (e.g., [17, p. 490, last paragraph] but, as the above analysis shows, that assumption is not necessary; orthonormality (paraunitariness in the uniform case) is really sufficient. Another subtle fact is that, when the quantizer noise enters an expander ($1n_i$ in Fig. 1(a)), it does not remain wide-sense stationary, but becomes cyclostationary [36]. This issue is correctly accommodated by the fact that we have averaged the output noise variance (over L samples) when obtaining (3.27).

In our derivations of the convolver coding gain, we assumed that the noise sources $q_k(n)$ are white, and uncorrelated with each other. Even these assumptions are not required in the subband coding case (i.e., the $g(n) = \delta(n)$ case). The orthonormality of the filter bank makes these assumptions unnecessary, as shown in [29] and [12, appendix sec. C.4.2]. On the other hand, it can be shown (Appendix C) that if the noise sources $q_k(n)$ are white and uncorrelated, then (3.43)–(3.45) can be obtained simply by assuming that the filters $f_k(n)$ have unit energy (that is, orthonormality is not necessary). Summarizing, the two sets of assumptions i) noise sources $q_k(n)$ are white and uncorrelated, and ii) filters are orthonormal are complementary to each other. Either one is sufficient to validate (3.43)–(3.45)!

IV. GENERAL ORTHOGONAL TRANSFORM CONVOLVER

The optimal coding gain (3.23) depends on the choice of the paraunitary matrix $E(z)$. A natural problem of interest here is the choice of optimal paraunitary $E(z)$ of a given degree J (for fixed number of channels M) which further maximizes the coding gain. In general this is a difficult problem, although some progress can be made in the special case where $J = 0$, i.e., $E(z)$ is a constant unitary matrix T . This is shown in Fig. 8(a). We will now consider the optimization problem for this special case. This special case is particularly attractive because the analysis filters $H_k(z)$ have length $\leq M$ (which could be much smaller than the lengths of $x(n)$ and $g(n)$). In this case the complexity of implementing the analysis and synthesis filters is negligible (compared to the complexity of the convolutions $x_k(n) * g_k^*(-n)$), and can therefore be disregarded. However, significant coding gain can still be achieved, as we will demonstrate.

With T taken to be unitary, i.e., $T^\dagger T = I$, the system is a paraunitary perfect reconstruction filter bank [2]. This is similar to the orthogonal transform coding system [17]. The convolution theorem (Theorem 2.1) clearly continues to hold in this case, and so do the coding gain expressions

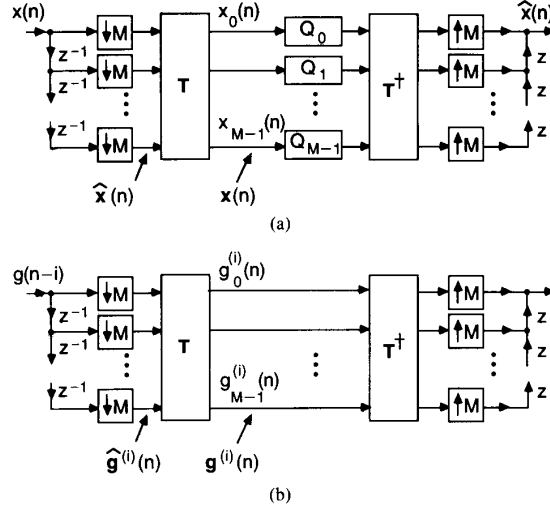


Fig. 8. The orthogonal transform convolver. (a) $x(n)$ is input to the filter bank, and (b) shifted $g(n)$ is input to the filter bank.

of the previous section. A special case of this system is the DFT filter bank, where T is the DFT or the IDFT matrix. We will now address the problem of finding the optimal T that maximizes the coding gain (3.17) under optimal bit allocation. It will again be assumed that the signals $x(n)$, $g(n)$ and the matrix T are real. We will first simplify the expression (3.17) by writing σ_{xk}^2 and α_k^2 directly in terms of T .

A. Expressions for σ_{xk}^2 and α_k^2

First refer to Fig. 8(a). Define the vectors $\hat{x}(n)$ and $x(n)$ as

$$\hat{x}(n) = \begin{bmatrix} x(Mn) \\ x(Mn - 1) \\ \vdots \\ x(Mn - M + 1) \end{bmatrix}, \quad x(n) = \begin{bmatrix} x_0(n) \\ x_1(n) \\ \vdots \\ x_{M-1}(n) \end{bmatrix}. \quad (4.1)$$

Then $x(n) = T\hat{x}(n)$. Assuming that $x(n)$ is WSS, the vector processes $\hat{x}(n)$ and $x(n)$ are WSS. Define the autocorrelations

$$\hat{R}_{xx} = E[\hat{x}(n)\hat{x}^\dagger(n)] \quad \text{and} \quad R_{xx} = E[x(n)x^\dagger(n)]. \quad (4.2)$$

Then

$$R_{xx} = T\hat{R}_{xx}T^\dagger. \quad (4.3)$$

The quantity σ_{xk}^2 is the diagonal element $[R_{xx}]_{kk}$ so that the product of these (which appears in the denominator of

(3.17)) is given by

$$\prod_{k=0}^{M-1} \sigma_{xk}^2 = \prod_{k=0}^{M-1} (\mathbf{T} \hat{\mathbf{R}}_{xx} \mathbf{T}^\dagger)_{kk}. \quad (4.4)$$

Next refer to Fig. 8(b). Define the vectors $\hat{\mathbf{g}}^{(i)}(n)$ and $\mathbf{g}^{(i)}(n)$ as in (3.20). We then have $\mathbf{g}^{(i)}(n) = \mathbf{T} \hat{\mathbf{g}}^{(i)}(n)$. Thus

$$\begin{aligned} \alpha_k^2 &= \sum_{i=0}^{M-1} \sum_n (\mathbf{g}^{(i)}(n) [\mathbf{g}^{(i)}(n)]^\dagger)_{kk} \\ &\quad (\text{from the definition (3.8)}) \\ &= \left(\mathbf{T} \sum_{i=0}^{M-1} \sum_n \hat{\mathbf{g}}^{(i)}(n) [\hat{\mathbf{g}}^{(i)}(n)]^\dagger \mathbf{T}^\dagger \right)_{kk} \\ &= (\mathbf{T} \hat{\mathbf{R}}_{gg} \mathbf{T}^\dagger)_{kk} \end{aligned} \quad (4.5)$$

where

$$\hat{\mathbf{R}}_{gg} = \sum_{i=0}^{M-1} \sum_n \hat{\mathbf{g}}^{(i)}(n) [\hat{\mathbf{g}}^{(i)}(n)]^\dagger. \quad (4.6)$$

Summarizing, the convolutional coding gain (3.17) can be expressed as

$$G_{TC}(M) = \frac{\sigma_x^2 \sum_n |g(n)|^2}{\left(\prod_{k=0}^{M-1} (\mathbf{T} \hat{\mathbf{R}}_{xx} \mathbf{T}^\dagger)_{kk} \prod_{k=0}^{M-1} (\mathbf{T} \hat{\mathbf{R}}_{gg} \mathbf{T}^\dagger)_{kk} \right)^{1/M}}. \quad (4.7)$$

The subscript TC stands for transform coding. The expression (4.7) holds under the optimal bit allocation condition (3.16). The unitary matrix \mathbf{T} should be chosen so as to minimize the product in the denominator.

B. Properties of the Matrices $\hat{\mathbf{R}}_{xx}$ and $\hat{\mathbf{R}}_{gg}$

The $M \times M$ matrix $\hat{\mathbf{R}}_{xx}$ is the autocorrelation matrix derived from a scalar WSS process $x(n)$, and is therefore Hermitian, Toeplitz, and positive semidefinite. It is also positive definite unless $x(n)$ is harmonic (i.e., the power spectrum is made of impulses $\delta(\omega - \omega_k)$). It can be shown that $\hat{\mathbf{R}}_{gg}$ also has all these properties, i.e., the Hermitian, Toeplitz, and positive definite unless $G(e^{j\omega})$ is made of impulses. (See Appendix D.) In fact it turns out that $[\hat{\mathbf{R}}_{gg}]_{km} = \sum_n g(n) g^*(n + k - m)$ so that it is a deterministic autocorrelation matrix.

The problem of finding the optimal transformation \mathbf{T} therefore reduces to the following: given the $M \times M$ Hermitian, Toeplitz and positive definite matrices $\hat{\mathbf{R}}_{xx}$ and $\hat{\mathbf{R}}_{gg}$, find a unitary matrix \mathbf{T} such that

$$\prod_{k=0}^{M-1} (\mathbf{T} \hat{\mathbf{R}}_{xx} \mathbf{T}^\dagger)_{kk} \prod_{k=0}^{M-1} (\mathbf{T} \hat{\mathbf{R}}_{gg} \mathbf{T}^\dagger)_{kk} \quad (4.8)$$

is minimized.

Given a Hermitian positive definite matrix \mathbf{P} , consider the product $\prod_{k=0}^{M-1} (\mathbf{T} \mathbf{P} \mathbf{T}^\dagger)_{kk}$ where \mathbf{T} is constrained to be unitary. It is known that this product is minimized if and only if the columns of \mathbf{T}^\dagger are eigenvectors of \mathbf{P} . (This is

how the traditional Karhunen-Loeve transform (KLT) is obtained [17]). Under this condition $\mathbf{T} \mathbf{P} \mathbf{T}^\dagger$ is diagonal. However, in our case, two positive definite matrices are involved. The problem of finding a single unitary matrix \mathbf{T} that minimizes the product (4.8) does not appear to have a simple, known solution.

If the matrices $\hat{\mathbf{R}}_{xx}$ and $\hat{\mathbf{R}}_{gg}$ are diagonalizable by the same unitary matrix \mathbf{T} , then this \mathbf{T} maximizes the coding gain. This condition for simultaneous diagonalization is equivalent to either of the following two conditions [37]:

- 1) $\hat{\mathbf{R}}_{xx}$ and $\hat{\mathbf{R}}_{gg}$ commute, i.e., $\hat{\mathbf{R}}_{xx} \hat{\mathbf{R}}_{gg} = \hat{\mathbf{R}}_{gg} \hat{\mathbf{R}}_{xx}$.
- 2) $\hat{\mathbf{R}}_{xx} \hat{\mathbf{R}}_{gg}$ is Hermitian.

For the special case of 2×2 real matrices (i.e., $M = 2$, and $x(n)$, $g(n)$, and \mathbf{T} are real), the above conditions are satisfied for the following reason: The matrices $\hat{\mathbf{R}}_{xx}$ and $\hat{\mathbf{R}}_{gg}$ are 2×2 symmetric Toeplitz, so that they are also circulant. But circulant matrices commute [38]. The two matrices are simultaneously diagonalizable by the unitary matrix

$$\mathbf{T} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}. \quad (4.9)$$

With this choice of \mathbf{T} the coding gain reduces to

$$G_{TC}(2) = \frac{1}{\sqrt{(1 - \rho_x^2)(1 - \rho_g^2)}} \quad (4.10)$$

where $\rho_x = E[x(n)x^*(n-1)]/\sigma_x^2$, and $\rho_g = \sum_n g(n)g^*(n-1)/\sum_n |g(n)|^2$. For example, if $\rho_x = \rho_g = 0.95$ then the coding gain is $G_{TC}(2) = 10.26$.

C. Bound on the Coding Gain

For a Hermitian positive definite matrix \mathbf{P} , we have $\prod_{i=0}^{M-1} [\mathbf{P}]_{ii} \geq \det \mathbf{P}$ with equality if and only if \mathbf{P} is diagonal. Using this we see that the gain (4.7) is bounded as

$$G_{TC}(M) \leq \frac{\sigma_x^2 \sum_n |g(n)|^2}{([\det \hat{\mathbf{R}}_{xx}][\det \hat{\mathbf{R}}_{gg}]^{1/M})^{1/M}} = \frac{\sigma_x^2 \sum_n |g(n)|^2}{(\det [\hat{\mathbf{R}}_{xx} \hat{\mathbf{R}}_{gg}])^{1/M}}. \quad (4.11)$$

V. NUMERICAL EXAMPLES

In the following examples, we will demonstrate the coding gains of the paraunitary convolvers. The signals $x(n)$, $g(n)$, and the number of subbands M are chosen as follows:

- 1) Number of subbands $M = 6$ in all cases.
- 2) Many choices of $g(n)$ are used, but all of these are such that $G(e^{j\omega})$ is low pass as demonstrated in Fig. 9. All choices have the same band edges. To obtain different stopband attenuations, we change the length of $g(n)$, but retain the same bandedges for $G(e^{j\omega})$.
- 3) The input signal $x(n)$ is taken to be an autoregressive process of order five (i.e., an AR(5) process). The autocorrelation coefficients $R(k)$, for $0 \leq k \leq 5$, are obtained from [17, table 2.2] (low-pass speech source).

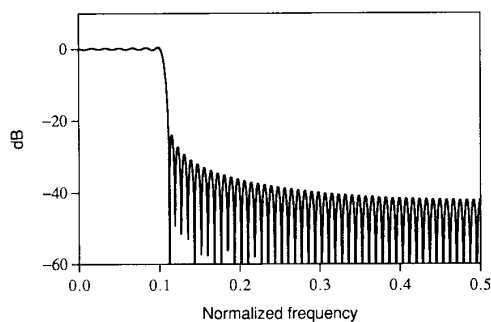


Fig. 9. A typical magnitude response of the filter $g(n)$ used in the experiment.

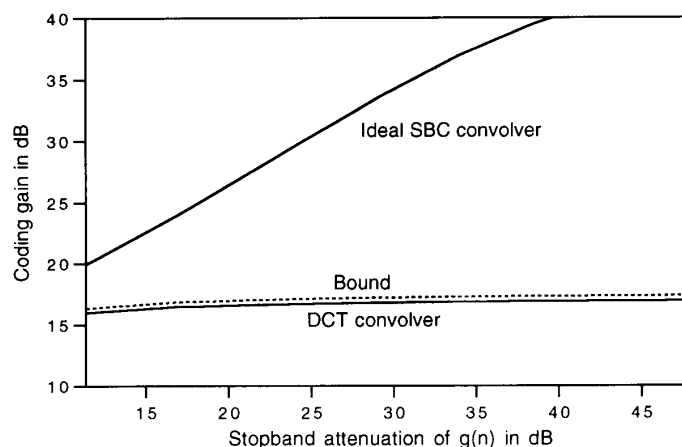


Fig. 10. Demonstration of the coding gains of paraunitary convolvers.

Where necessary, the power spectrum $S_{xx}(e^{j\omega})$ is computed as $S_{xx}(e^{j\omega}) = \alpha / |1 + \sum_{n=1}^5 a_n e^{-j\omega n}|^2$ where a_n are the autoregressive coefficients (obtainable by solving the optimal fifth-order linear-prediction problem [17]).

Fig. 10 shows the coding gain of the paraunitary convolver (with optimal bit allocation) as a function of the stopband attenuation of $G(e^{j\omega})$, for three cases. The top-most curve corresponds to the ideal SBC convolver. In other words, the analysis and synthesis filters are as in Fig. 6 (ideal brick-wall filters (3.24)). The bottom curve is for the DCT convolver, that is an orthogonal transform convolver (Fig. 8) in which the matrix T is taken to be the 6×6 DCT matrix. (Four types of DCT matrix have been defined in the literature; we have used the one in [17, (12.157)].)⁵ The middle curve shows the upper bound (4.11) for the orthogonal transform convolver. It is interesting to note that the DCT system is only about 0.5 dB worse than the bound. The ideal brick wall SBC convolver is significantly better than the DCT convolver. The DCT convolver, however, is very simple to implement (much less expensive than good filters approximating the

ideal SBC filters). In all the above cases the coding gain improves with the attenuation of $G(e^{j\omega})$ because the AM/GM ratio in (3.17) improves.

In the above experiment suppose we take $g(n) = \delta(n)$. Then the coding gain of the convolver is equal to the coding gain of the traditional subband coding system. For the ideal SBC filters, this value is $G = 6.72$ dB, and for transform coding with DCT this is 5.3 dB (consistent with experiments on speed coding; for example, see [17, p. 542]). Thus, the large additional gain seen in Fig. 10 is contributed by the filter $G(e^{j\omega})$ participating in the subband convolver.

We have not shown plots of the coding gain with respect to the number of channels M , as it does not reveal more insights than what is already known in subband coding practice [17], [18], [39].

VI. BIORTHONORMAL FILTER-BANK CONVOLVERS

While this paper was being reviewed, it was pointed out by S.-M. Phoong (graduate student, California Institute of Technology) that a previous paper by this author had an example of a "filter-bank convolution theorem" in hidden form (see the figures in [2, p. 84]). Furthermore, while it assumed the filter bank to have perfect re-

⁵The motivation for the use of the DCT is that in traditional speech coding, it is known to be an excellent substitute for the optimal (KLT) transform.

construction, the example worked just fine even without orthonormality!

This motivated the author to generalize the results of Section II for arbitrary filter banks: Consider again Fig. 1(a), where $H_k(z)$ and $F_k(z)$ denote the analysis and synthesis filters, respectively, of an M channel filter bank with subband decimation ratios n_k . Assume maximal decimation, i.e., $\sum_k 1/n_k = 1$. It can be shown (Appendix B) that the perfect reconstruction property (i.e., $\hat{x}(n) = x(n)$ in Fig. 1(a)) is ensured by the condition

$$(F_k(z)H_m(z))_{l_{n_k, m}} = \delta(k - m) \quad (\text{biorthonormality}) \quad (6.1)$$

where $n_{k, m} = \text{gcd}(n_k, n_m)$. The above property, called biorthonormality, reduces to (2.18) in the orthonormal case because of the condition $F_k(z) = \tilde{H}_k(z)$. We will now prove the following result.

Theorem 6.1: Biorthonormal convolution theorem. Let the two signals $x(n)$ and $g(n)$ be passed through the two different analysis filter banks as shown in Fig. 11. Assume that $\{H_k(z)\}$ and $\{F_k(z)\}$ are, respectively, the analysis and synthesis filters of a maximally decimated biorthonormal filter bank with decimation ratios n_k . Then the convolution $x(n) * g(n)$ is related to the convolutions of the subband signals $x_k(n)$ and $g_k(n)$ as follows

$$(x(n) * g(n))_{l_L} = \sum_{k=0}^{M-1} (x_k(n) * g_k(n))_{l_{p_k}} \quad (6.2)$$

where all notations are as in Section II-B, that is, $L = \text{lcm}\{n_k\}$ and $L = n_k p_k$. \diamond

Proof: If the subband signals $x_k(n)$ are passed through the synthesis bank (as in Fig. 1(a) with quantizers ignored), then the perfect reconstruction property ensures

$$X(z) = \sum_{k=0}^{M-1} X_k(z^{n_k}) F_k(z). \quad (6.3)$$

Now, if we interchange each analysis filter $H_k(z)$ with the corresponding synthesis filter, the perfect reconstruction property is not affected (since (6.1) remains valid). Thus, if the subband signals $g_k(n)$ in Fig. 11(b) are passed through a synthesis bank with filters $H_k(z)$, we get back $g(n)$. That is,

$$G(z) = \sum_{m=0}^{M-1} G_m(z^{n_m}) H_m(z). \quad (6.4)$$

Computing $X(z)G(z)$ from these, and proceeding as in Theorem 2.2, we get

$$(X(z)G(z))_{l_L} = \sum_{k=0}^{M-1} (X_k(z)G_k(z))_{l_{p_k}} \quad (6.5)$$

which proves (6.2). $\nabla\nabla\nabla$

Comments

1) Thus, we can convolve the two signals by decoupled convolutions in the subbands, provided we use the anal-

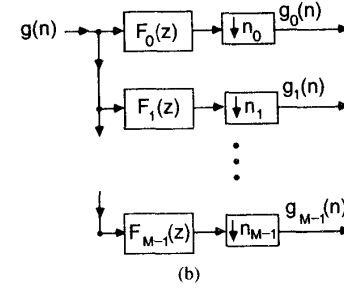
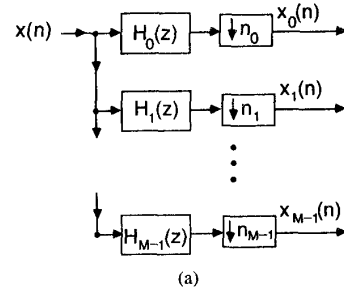


Fig. 11. Pertaining to the convolution theorem for biorthonormal filter banks. (a) Decomposition of $x(n)$, and (b) decomposition of $g(n)$.

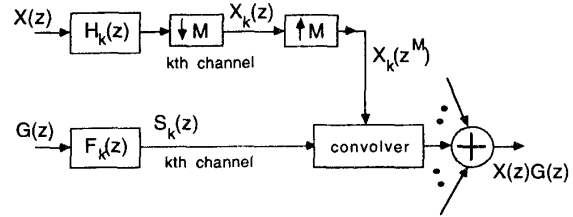


Fig. 12. Implementation of the biorthonormal convolver.

ysis bank to decompose one signal $x(n)$ and the synthesis bank to decompose the other signal.

2) *Obtaining all samples.* To obtain all the samples of $x(n) * g(n)$, we have to repeat (6.2) with shifted versions of $g(n - i)$. In practice, this can be done indirectly. For example, consider the uniform case, where the orthonormal system was implemented as in Fig. 3. The modified implementation for the biorthonormal case is shown in Fig. 12, where the box labelled "convolver" computes $S_k(z)X_k(z^M)$. For the nonuniform case, we can obtain an implementation by first converting the filter bank to a uniform system as described in Appendix B-3.

3) Note that, unlike in Theorems 2.1 and 2.2, we directly obtain the convolution $x(n) * g(n)$ rather than the correlation $x(n) * g^*(-n)$.

4) For the special case of orthonormal filter banks, $F_k(z) = \tilde{H}_k(z)$. So the two signals $x(n)$ and $g(n)$ are decomposed by $\{H_k(z)\}$ and $\{\tilde{H}_k(z)\}$, respectively. This is similar to decomposing $x(n)$ and $g^*(-n)$ with the same analysis bank $H_k(z)$, so that Theorem 2.2 is obtained as a special case.

5) *Coding gain.* For the biorthonormal case, we have omitted the derivation of bit allocation and coding gain

formulas. The derivation requires some modifications of the previous sections. Appendix C will be of some help in the derivation.

VII. CONCLUDING REMARKS

In this paper we have introduced convolution theorems for filter bank transformers. Both uniform and nonuniform decimation ratios were considered, and orthonormal as well as biorthonormal cases were addressed. All the theorems are such that the original convolution reduces to a sum of shorter, decoupled, convolutions in the subbands. That is, there is no need to have cross convolution between subbands.

For the orthonormal case, expressions for optimal bit allocation and the optimized coding gain were derived. The contribution to coding gain comes partly from the nonuniformity of the signal spectrum $S_{xx}(e^{j\omega})$, and partly from nonuniformity of the filter spectrum $|G(e^{j\omega})|^2$. With $g(n)$ taken to be the unit pulse function $\delta(n)$, the coding gain expressions reduce to those for traditional subband and transform coding, many of which are well known.

The filter-bank convolver has about the same computational complexity as a traditional convolver, if the analysis bank has small complexity compared to the convolution itself. Such, indeed, is the situation in the special case of the orthogonal transform convolver (Fig. 8) where the analysis filter bank has filter lengths $\leq M$ (number of bands). In spite of this simplicity, the coding gain obtainable can already be quite significant. Even though there is no closed form expression for the optimal orthogonal convolver matrix T , we could derive an upper bound for this (for fixed M), and the DCT matrix offers a gain very close to this bound for the case of speech signals.

A. Putting Things in Perspective

The power of orthogonality in the reduction of double summations into single ones has been used over and over again, in all fields of science and engineering. And yet, only very special cases of convolution theorems have been reported in the past. To explain the reason for this, let us switch for a moment into heuristic mode, and imagine that the samples of two sequences $x_1(n)$ and $x_2(n-i)$ (where i is a fixed shift index) are collected into vectors x_1 and x_2 . Let T be a unitary transformation (i.e., $T^\dagger T = I$), and let $y_1 = Tx_1$ and $y_2 = Tx_2$. Then the unitariness implies $y_2^\dagger y_1 = x_2^\dagger x_1$. (This is essentially Parseval's relation.) Starting from this and varying i , one could obtain "convolution type of theorems." The reason why this is not as simple as it looks is due to the sizes of the vectors and matrices. If all of these are infinite dimensional then the result is of little use. In the finite size case, we have to account for the fact that the size grows after convolution (or use circular convolutions; see below). So the "border effects" are crucial. In fact, well-known convolution theorems differ from each other primarily in the way they handle this issue.

The most well-known (perhaps earliest) successful con-

volution theorems were based on the traditional continuous and discrete time Fourier transforms. Then came the circular convolution methods [13], which work for finite length sequences (which can be imagined periodic). They can be nicely adopted to perform finite linear convolutions (or infinite length ones, by sectioning). The circular convolution theorems, however, hold only for special types of orthogonal transforms, with the "primitive root property" (see [40, sec. I]). Examples are the DFT and the number theoretic transforms [40].

On the other hand, for certain types of orthogonal transforms, such as the discrete cosine and sine transform (DCT, DST), the convolution theorems are more complicated. See [41]–[44], and references at the bottom of [43, p. 459]. In these situations, one starts with a finite length sequence, and constructs a symmetric or antisymmetric sequence (nearly two times longer) to which the transform is applied. The details depend on the type of DCT and so forth (there are at least four known types). For more arbitrary orthogonal transforms, it appears that convolution theorems have not been reported earlier.

It seems, therefore, that the really nontrivial issue in any kind of convolution theorem has to do with the fact that we wish to use finite transforms (computable in finite time), and need to take care of border effects one way or the other. The details of this depend on the coefficients of the transform matrix (DFT, DCT, etc.) and the type of convolution (linear, circular, etc.).

In all these earlier techniques, the attempt is always to convert "convolution" in one domain into point by point multiplication in the other domain. The result presented in this paper, however, differs in this respect. Thus, once we pass into the subband domain, we perform "subband by subband convolution." If we view the filter bank as a transformer from time to time-frequency, then the transform domain quantities are $x_k(m)$ (Fig. 1(a)) where (m, k) is the time-frequency index. We perform convolution with respect to "time" m and add up the results for all "frequencies" k . In other words, we do not perform point by point multiplication in the (m, k) domain. This is why the theorems work for any type of inputs (infinite or finite); and all convolutions are infinite length, linear convolutions. As we have shown, these results work for all invertible filter banks—orthonormal, biorthonormal, nonuniform, and so forth. Further details of the filter-bank coefficients have essentially no role.

B. Generalizations, and Open Problems

The results of this paper naturally lead us to ask if the convolution theorems are true for other types of filter banks, e.g., those with rational decimation ratios. Using a vector space approach, the results have been generalized [45], and hold even for multidimensional filter banks with arbitrary, nonuniform, rational decimation matrices.

Some issues still need to be addressed. For example, in the 1D orthogonal convolver of Section IV, it is still of some theoretical interest to find the best unitary T that

minimizes (4.8) (i.e., maximizes the coding gain). A second problem is the derivation of the bit allocation and coding gain formulas for the biorthonormal case. A third problem is the application of these ideas to subband adaptive filtering [46], [47]. In these techniques one usually has to allow "cross terms" between subbands; it might be worth trying to reconfigure the adaptive filtering system so that the decoupling of the subbands can somehow be exploited.

A fourth problem is the extension of the results to the case of continuous time orthonormal wavelet transforms. It is well known [7] that a large class of signal $x_a(t)$ can be expressed in the form

$$x_a(t) = \sum_{k,m} X_{DWT}(k, m) \psi_{k,m}(t) \quad (7.1)$$

where $\psi_{k,m}(t)$ are a class of orthonormal basis functions derived from a wavelet function $\psi(t)$ by dilations and shifts:

$$\psi_{k,m}(t) = 2^{-k/2} \psi(2^{-k}t - m). \quad (7.2)$$

Here (m, k) can be regarded as the transform domain (time-frequency domain). Notice that, unlike the filter bank case, we have an infinite number of values of k here. It will be interesting to find convolution theorems for this kind of orthonormal decompositions.

APPENDIX A

EQUIVALENCE OF (2.15) AND (2.17)

If $k = m$ we can rewrite (2.15) as

$$\sum_n f_k(n) f_k^*(n + n_k(l - i)) = \delta(l - i)$$

and (2.17) as

$$\sum_n f_k(n) f_k^*(n + n_k p) = \delta(p).$$

Evidently these imply each other.

Next let $k \neq m$. First assume that (2.15) holds. Recall $n_{k,m} = \gcd(n_k, n_m)$. Thus, there exist integers a and b such that $n_k a - n_m b = n_{k,m}$. Therefore, given any integer p there exists integers l and i such that $n_k l - n_m i = n_{k,m} p$. Thus the left-hand side in (2.17) can always be rewritten to resemble the left-hand side of (2.15). Since $k \neq m$, this left-hand side is indeed zero, so that the left-hand side of (2.17) is zero as well. Conversely, let (2.17) be true. Given a pair of integers l, i we can always write $n_k l - n_m i = n_{k,m} p$ for some integer p . So the left side of (2.15) can be rewritten to resemble the left side of (2.17). Since $k \neq m$, (2.17) says that this is zero, so that the same follows for (2.15).

APPENDIX B

NONUNIFORM FILTER BANKS

1. Biorthonormality and Perfect Reconstruction

Consider the nonuniform maximally decimated system of Fig. 1(a). Suppose the synthesis filters have the potentiality for perfect reconstruction (i.e., there exist analysis

filters such that these synthesis filters give perfect reconstruction). This means (similar to (2.12b))

$$x(n) = \sum_{k=0}^{M-1} \sum_l \alpha_k(l) f_k(n - n_k l) \quad (B.1)$$

for any $x(n)$, for appropriate choices of $\alpha_k(l)$. Thus, in order to have perfect reconstruction, it is necessary for the set of basis functions $\{\xi_{k,l}(n)\} \triangleq \{f_k(n - n_k l)\}$ to be "complete." We shall assume that this is the case.

We now show that if the analysis filters $H_m(z)$ are chosen to satisfy the following condition:

$$\begin{aligned} \sum_n f_m(n - n_m i) h_k(-n + n_k l) \\ = \delta(k - m) \delta(l - i) \end{aligned} \quad (\text{biorthonormality}) \quad (B.2)$$

then the filter bank has perfect reconstruction. Note that the left side of the above equation, viewed as a function of l , is the n_k -fold decimated version of the convolution $f_m(n - n_m i) * h_k(n)$. Thus, taking z transforms, (B.2) is equivalent to

$$(z^{-inm} F_m(z) H_k(z))_{\downarrow n_k} = z^{-i} \delta(k - m). \quad (B.3)$$

Suppose we choose the filter bank input to be $x(n) = f_m(n - n_m i)$, that is $X(z) = z^{-inm} F_m(z)$ for some m , and some i . Then (B.3) implies that the decimated output of the analysis filter $H_k(z)$ is zero, for $k \neq m$. And the decimated output of $H_m(z)$ has the z -transform z^{-i} . If this set of signals is passed through the synthesis bank, the reconstructed output has z -transform $z^{-inm} F_m(z)$. That is $\hat{X}(z) = z^{-inm} F_m(z) = X(z)$ indeed. Since the filter bank is linear (though time varying), we conclude that any input $x(n)$ of the form (B.1) is perfectly recovered, i.e., $\hat{x}(n) = x(n)$.

As in Section II-B, let $n_{k,m} = \gcd(n_k, n_m)$. Then we can rewrite the condition (B.2) to obtain

$$(F_m(z) H_k(z))_{\downarrow n_{k,m}} = \delta(k - m) \quad (\text{biorthonormality}). \quad (B.4)$$

The proof is similar to the one in Appendix A.

Summarizing, the biorthonormality condition (B.4) ensures perfect reconstruction for any $x(n)$. It is also clear from (B.4) that if we interchange each $H_k(z)$ with the corresponding $F_k(z)$, the perfect reconstruction property is preserved.

2. Orthonormality

Recall the definition of orthonormality from Section II-B. If the synthesis filters are orthonormal (i.e., $(F_m(z) \tilde{F}_k(z))_{\downarrow n_{k,m}} = \delta(k - m)$) then we see from (B.4) that the choice of analysis filters $H_m(z) = \tilde{F}_m(z)$ gives perfect reconstruction. We will now prove the generalized power complementary property (2.19). For this, note that $\hat{X}(z)$ in Fig. 1(a) (ignoring quantizers) can always be expressed as

$$\hat{X}(z) = \sum_{k=0}^{M-1} F_k(z) \frac{1}{n_k} \sum_{m=0}^{n_k-1} H_k(z W_{n_k}^m) X(z W_{n_k}^m) \quad (B.5)$$

where $W_{n_k} = e^{-j2\pi/n_k}$. A perfect reconstruction system is, in particular alias-free (i.e., terms with $m \neq 0$ vanish) so that this gives

$$\frac{\hat{X}(z)}{X(z)} = \sum_{k=0}^{M-1} \frac{F_k(z)H_k(z)}{n_k}. \quad (\text{B.6})$$

Perfect reconstruction implies that this is unity. Substituting $H_k(z) = \tilde{F}_k(z)$, we immediately obtain (2.19).

3. Equivalent Uniform Filter Bank

Let $L = \text{lcm}\{n_k\}$, and $L = n_k p_k$ as usual. It can be shown [21, Fig. 3] that the M -channel nonuniform filter bank of Fig. 1(a) can be redrawn as an L channel maximally decimated uniform system (i.e., with equal decimation ratios L in all channels). The M sets of analysis and synthesis filters $\{H_k(z), F_k(z)\}$, of the nonuniform system are replaced with the L sets of filters $\{H'_k(z), F'_k(z)\}$ in the uniform system. Since $L > M$ unless all n_k are identical, we say that $\{H'_k(z), F'_k(z)\}$ is the “bigger system.” This equivalence can sometimes be exploited to derive useful conclusions [21], [22].

To obtain the equivalent uniform system, consider the k th channel shown separately in Fig. 13(a). It can be redrawn as shown in Fig. 13(b). This follows from the fact that a p_k -channel uniform filter bank with analysis filters z^{-i} and synthesis filters z^i , $0 \leq i \leq p_k - 1$, has perfect reconstruction. Fig. 13(b) can further be redrawn as in (c) by the use of noble identities [2], and by using $L = n_k p_k$. Thus, the k th channel can be expanded into p_k channels. Altogether we therefore have $\sum_k p_k = L$ channels with decimation ratio L in each channel. An integer k in $0 \leq k \leq M - 1$ and an integer i in $0 \leq i \leq p_k - 1$ uniquely identify the analysis and synthesis filters of the uniform system as $z^{-in_k} H_k(z)$ and $z^{in_k} F_k(z)$. The nonuniform system is denoted by $\{H_k(z), F_k(z), n_k\}$ where $0 \leq k \leq M - 1$, and the uniform system by $\{H'_k(z), F'_k(z), L\}$, where $0 \leq k \leq L - 1$. The equivalence of the two systems means that, for a given input $x(n)$, the output $\hat{x}(n)$ is the same for the two filter banks.

Fact B.1: Consider the following properties of the nonuniform system $\{H_k(z), F_k(z), n_k\}$.

- 1) $\hat{x}(n) = x(n)$ (perfect reconstruction).
- 2) $F_k(z) = \tilde{H}_k(z)$, i.e., $f_k(n) = h_k^*(-n)$ (time reversal property).
- 3) $(F_k(z)H_m(z))_{in_k, m} = \delta(k - m)$ (biorthonormality).
- 4) $(F_k(z)\tilde{F}_m(z))_{in_k, m} = \delta(k - m)$ (orthonormality).

If any one of these properties is true, then the corresponding property holds for the uniform system $\{H'_k(z), F'_k(z), L\}$. The converse is also true. \diamond

Proof: Proofs are required only for properties 3 and 4. For this consider the set of sequences

$$\begin{aligned} \{\xi_{m,i}(n)\} &= \{f_m(n - n_m i)\}, \\ 0 \leq m \leq M - 1, \quad -\infty \leq i \leq \infty \end{aligned} \quad (\text{B.7})$$

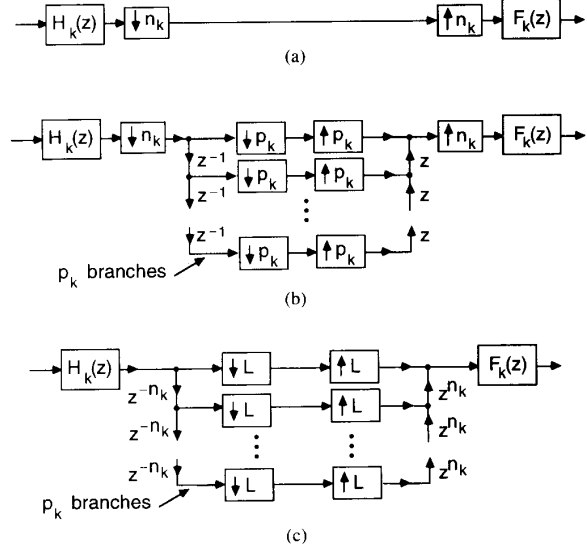


Fig. 13. Redrawing a nonuniform filter bank as a uniform filter bank. (a) The k th channel, (b) the redrawn version, and (c) expanded version of the k th channel.

appearing in (B.2). For the uniform L -channel system, the corresponding set of sequences are

$$\begin{aligned} \{f_m(n + jn_m - Li)\}, \quad 0 \leq j \leq p_m - 1, \\ 0 \leq m \leq M - 1, \quad -\infty \leq i \leq \infty. \end{aligned} \quad (\text{B.8})$$

Since $L = n_m p_m$, this set is the same as

$$\begin{aligned} \{f_m(n + n_m(j - ip_m))\}, \quad 0 \leq j \leq p_m - 1, \\ 0 \leq m \leq M - 1, \quad -\infty \leq i \leq \infty. \end{aligned} \quad (\text{B.9})$$

Clearly the two sets of sequences (B.7) and (B.9) are identical. A similar statement follows for the analysis filters as well. Consequently, the uniform system $\{H_k(z), F_k(z), n_k\}$ is biorthonormal if and only if the uniform system $\{H'_k(z), F'_k(z), L\}$ is biorthonormal. Identical reasoning can be given for orthonormality. $\nabla\nabla\nabla$

Fact B.2: For a nonuniform maximally decimated system with analysis filters $H_k(z)$, synthesis filters $F_k(z)$, and decimation ratios n_k , consider the following three properties:

- 1) $\hat{x}(n) = x(n)$ (perfect reconstruction).
- 2) $F_k(z) = \tilde{H}_k(z)$, i.e., $f_k(n) = h_k^*(-n)$ (time reversal property).
- 3) $(F_k(z)\tilde{F}_m(z))_{in_k, m} = \delta(k - m)$ (orthonormality).

If any two of these is true, then the remaining property is also true. \diamond

Proof: For the uniform case (i.e., $n_k = M$ for all k) this has been proved in [12] (Theorem 6.2.1). For the nonuniform case this follows by defining the uniform L channel system as above, and invoking Fact B.1. \diamond

APPENDIX C

ON THE WHITE, UNCORRELATED ASSUMPTION

We will show that if the quantizer noise source $q_k(n)$ due to k th subband quantizer (Fig. 4) is white, and if the noise sources are uncorrelated for different values of k ,

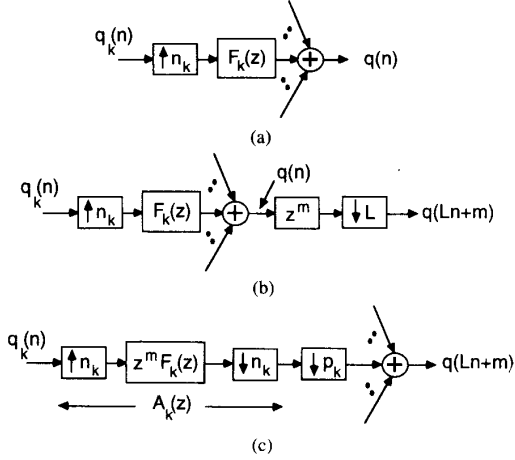


Fig. 14. (a) The k th branch of the synthesis bank with a noise source at its input, (b) insertion of a delay and a decimator, and (c) rearrangement.

then the expression (3.43) for output noise variance holds, as long as $F_k(z)$ have unit energy. Fig. 14(a) shows the k th synthesis filter branch, with $q_k(n)$ as its input.

In the presence of the expanders (denoted $\uparrow n_k$ in the figures) the total output noise $q(n)$ is not wide-sense stationary (WSS), unless the filters $F_k(z)$ are ideal [36]. In order to handle this situation, imagine that the function z^m and a decimator L are inserted at the output as shown in Fig. 14(b). Using $L = n_k p_k$ we can redraw this as in Fig. 14(c). Then the system indicated $A_k(z)$ in the figure is a linear time invariant system with impulse response $f_k(n_k n + m)$, i.e., it is a polyphase component of $F_k(z)$. (See "polyphase identity" on [12, p. 133].) Since the decimators ($\downarrow p_k$) retain the WSS property of random processes, it is clear that the decimated output noise $q(Ln + m)$ is WSS. Using standard techniques, and the assumptions that the sources $q_k(n)$ are jointly WSS, zero-mean white, and uncorrelated for different values of k , the variance of $q(Ln + m)$ is

$$E[|q(Ln + m)|^2] = \sum_{k=0}^{M-1} \sigma_{q_k}^2 \sum_n |f_k(n_k n + m)|^2 \quad (C.1)$$

where $\sigma_{q_k}^2$ is the variance of $q_k(n)$. We have to average this for $0 \leq m \leq L - 1$, to obtain a constant answer. Thus,

$$\sigma_q^2 = \text{Average output noise variance} = \frac{1}{L} \sum_{k=0}^{M-1} \sum_{m=0}^{L-1} \sigma_{q_k}^2 \cdot \sum_n |f_k(n_k n + m)|^2. \quad (C.2)$$

Since $L = p_k n_k$, this simplifies to

$$\sigma_q^2 = \frac{1}{L} \sum_{k=0}^{M-1} \sigma_{q_k}^2 p_k \sum_n |f_k(n)|^2 = \sum_{k=0}^{M-1} \frac{\sigma_{q_k}^2}{n_k} \sum_n |f_k(n)|^2. \quad (C.3)$$

This indeed reduces to

$$\sigma_q^2 = \sum_{k=0}^{M-1} \frac{\sigma_{q_k}^2}{n_k} \quad (C.4)$$

when each filter $f_k(n)$ has unit energy.

Summarizing, the expressions (3.43)–(3.45) for the subband coder are true under the assumptions that i) the quantizer noise sources $q_k(n)$ are white, and uncorrelated for different k , and ii) the synthesis filters $F_k(z)$ have unit energy.

APPENDIX D

NONSINGULARITY OF \hat{R}_{xx} AND \hat{R}_{gg}

Since \hat{R}_{xx} is the autocorrelation matrix obtained from a scalar WSS process $x(n)$, it is positive semidefinite. It is therefore positive definite if and only if it is nonsingular. If this matrix is singular, then there exists $\mathbf{v} \neq \mathbf{0}$ such that $\mathbf{v}^T \hat{R}_{xx} \mathbf{v} = 0$, i.e., $E[|\mathbf{v}^T \hat{\mathbf{x}}(n)|^2] = 0$, i.e., $\mathbf{v}^T \hat{\mathbf{x}}(n) = 0$. In other words, there exists an FIR filter $V(z) \triangleq v_0^* + v_1^* z^{-1} + \dots + v_{M-1}^* z^{-(M-1)}$ such that the output in response to the WSS process $x(n)$ is zero. Thus if $S_{xx}(e^{j\omega})$ denotes the power spectrum of $x(n)$, then the power spectrum of the output is $S_{xx}(e^{j\omega}) |V(e^{j\omega})|^2 = 0$. Since the FIR filter $V(z)$ can have at most $M - 1$ zeros on the unit circle, this means that the power spectrum has the form $S_{xx}(e^{j\omega}) = \sum_{k=1}^{M-1} c_k \delta(\omega - \omega_k)$, i.e., $x(n)$ is a harmonic process. Thus, unless $x(n)$ is harmonic, \hat{R}_{xx} is positive definite. This is a well-known fact [48], and is reviewed here only for completeness.

Next consider \hat{R}_{gg} defined in (4.6). Using the definition of $\hat{\mathbf{g}}^{(i)}(n)$ in (3.20) we see that

$$\begin{aligned} [\hat{R}_{gg}]_{pq} &= \sum_{i=0}^{M-1} \sum_n g(Mn - i - p) g^*(Mn - i - q) \\ &= \sum_l g(l - p) g^*(l - q) = R_{gg}(q - p) \end{aligned} \quad (D.1)$$

where $R_{gg}(k)$ is the deterministic autocorrelation of the sequence $g(n)$. Thus \hat{R}_{gg} is a deterministic autocorrelation matrix and has all the properties of \hat{R}_{xx} . It can be written as

$$\hat{R}_{gg} = \sum_n \begin{bmatrix} g(n) \\ g(n-1) \\ \vdots \\ g(n-M+1) \end{bmatrix} \cdot [g^*(n) \ g^*(n-1) \ \dots \ g^*(n-M+1)]. \quad (D.2)$$

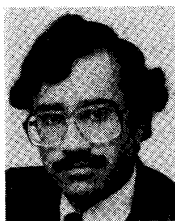
If this is singular, then there exists a vector $\mathbf{c} \neq \mathbf{0}$ such that $\mathbf{c}^T \hat{R}_{gg} \mathbf{c} = 0$. Thus, for each n in (D.2), we must have $c_0^* g(n) + c_1^* g(n-1) + \dots + c_{M-1}^* g(n-M+1) = 0$, where at least one c_i is nonzero. Proceeding as in the previous paragraph, we see that this happens only if $G(e^{j\omega})$ is either zero or made of at most $M - 1$ impulses.

ACKNOWLEDGMENT

The author thanks T. Chen and S.-M. Phoong, graduate students at Caltech, for interesting discussions and their valuable comments on the paper. S.-M. Phoong also generated the numerical examples presented in Section V. S. Martucci of the Georgia Institute of Technology sent the author some references on DCT convolutions, including his preprint on that topic. One of the reviewers drew his attention to a recent doctoral dissertation [49], which addresses a number of issues on subband coding. When this paper was on its way to the press, Prof. M. Vetterli of Columbia University drew his attention to a very interesting monograph [50], which addresses convolution using DFT filter banks.

REFERENCES

- [1] R. E. Crochiere and L. R. Rabiner, *Multirate Digital Signal Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1983.
- [2] P. P. Vaidyanathan, "Multirate digital filters, filter banks, polyphase networks, and applications: A tutorial," *Proc. IEEE*, vol. 78, pp. 56-93, Jan. 1990.
- [3] M. Vetterli, "A theory of multirate filter banks," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-35, pp. 356-372, Mar. 1987.
- [4] M. Vetterli, "Running FIR and IIR filtering using multirate filter banks," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-36, pp. 730-738, May 1988.
- [5] M. J. T. Smith and T. P. Barnwell, III, "Exact reconstruction techniques for tree-structured subband coders," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 34, pp. 434-441, June 1986.
- [6] I. Daubechies, "Orthonormal bases of compactly supported wavelets," *Commun. Pure Appl. Math.*, vol. 4, pp. 909-996, Nov. 1988.
- [7] I. Daubechies, "The wavelet transform, time-frequency localization, and signal analysis," *IEEE Trans. Inform. Theory*, vol. 36, pp. 961-1005, Sept. 1990.
- [8] S. Mallat, "A theory for multiresolution signal decomposition: The wavelet representation," *IEEE Trans. Patt. Anal. Machine Intell.*, vol. 11, pp. 674-693, July 1989.
- [9] A. H. Tewfik, D. Sinha, and P. E. Jorgensen, "On the optimal choice of a wavelet of signal representation," *IEEE Trans. Inform. Theory*, vol. 38, pp. 747-765, Mar. 1992.
- [10] O. Rioul and M. Vetterli, "Wavelets and signal processing," *IEEE Signal Processing Mag.*, pp. 14-38, Oct. 1991.
- [11] M. Vetterli and C. Herley, "Wavelets and filter banks: Theory and design," *IEEE Trans. Signal Processing*, vol. 40, no. 9, pp. 2207-2232, Sept. 1992.
- [12] P. P. Vaidyanathan, *Multirate Systems and Filter Banks*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [13] A. V. Oppenheim, A. S. Willsky, and I. T. Young, *Signals and Systems*. Englewood Cliffs, NJ: Prentice-Hall, 1983.
- [14] P. P. Vaidyanathan, "Theory and design of M -channel maximally decimated quadrature mirror filters with arbitrary M , having perfect reconstruction property," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-35, pp. 476-492, Apr. 1987.
- [15] P. P. Vaidyanathan, "Quadrature mirror filter banks, M -band extensions, and perfect reconstruction techniques," *IEEE ASSP Mag.*, vol. 4, pp. 4-20, July 1987.
- [16] R. E. Crochiere, "Subband coding," *Bell Syst. Tech. J.*, vol. 60, pp. 1633-1654, Sept. 1981.
- [17] N. S. Jayant and P. Noll, *Digital Coding of Waveforms*. Englewood Cliffs, NJ: Prentice-Hall, 1984.
- [18] J. W. Woods, *Subband Image Coding*. Kluwer, 1991.
- [19] Y. Huang and P. M. Schultheiss, "Block quantization of correlated Gaussian random variables," *IEEE Trans. Commun. Syst.*, pp. 289-296, Sept. 1963.
- [20] A. Segall, "Bit allocation and encoding for vector sources," *IEEE Trans. Inform. Theory*, pp. 162-169, Mar. 1976.
- [21] P.-H. Hoang and P. P. Vaidyanathan, "Nonuniform multirate filter banks: Theory and design," in *Proc. IEEE Int. Symp. Circuits Syst.*, Portland, OR, May 1989, pp. 371-374.
- [22] J. Kovačević and M. Vetterli, "Perfect reconstruction filter banks with rational sampling rate changes," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, Toronto, Canada, May 1991, pp. 1785-1788.
- [23] K. Nayebi, T. P. Barnwell, III, and M. J. T. Smith, "The design of perfect reconstruction nonuniform band filter banks," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, Toronto, Canada, May 1991, pp. 1781-1784.
- [24] A. Soman, and P. P. Vaidyanathan, "On orthonormal wavelets and paraunitary filter banks," *IEEE Trans. Signal Processing*, vol. 41, no. 3, pp. 1170-1184, Mar. 1993.
- [25] R. Coifmann, Y. Meyer, S. Quake, and V. Wickerhauser, "Signal processing with wavelet packets," Numer. Algorithm Res. Group, Yale Univ., 1990.
- [26] C. W. Barnes, B. N. Tran, and S. H. Leung, "On the statistics of fixed-point roundoff error," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-33, pp. 595-606, June 1985.
- [27] E. Beckenbach and R. Bellman, *An Introduction to Inequalities*. Random House, 1961.
- [28] H. S. Malvar, "Lapped transforms for efficient transform/subband coding," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 38, pp. 969-978, June 1990.
- [29] A. Soman and P. P. Vaidyanathan, "Coding gain in paraunitary analysis/synthesis systems," *IEEE Trans. Signal Processing*, vol. 41, no. 5, pp. 1824-1835, May 1993.
- [30] H. J. Nussbaumer, "Pseudo QMF filter bank," *IBM Tech. Disclosure Bull.*, vol. 24, pp. 3081-3087, Nov. 1981.
- [31] J. H. Rothweiler, "Polyphase quadrature filters, a new subband coding technique," in *Proc. IEEE Int. Conf. ASSP*, Boston, MA, Apr. 1983, pp. 1980-1983.
- [32] R. V. Cox, "The design of uniformly and nonuniformly spaced pseudo QMF," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-34, pp. 1090-1096, Oct. 1986.
- [33] J. Masson and Z. Picel, "Flexible design of computationally efficient nearly perfect QMF filter banks," in *Proc. IEEE Int. Conf. ASSP*, Tampa, FL, Mar. 1985, pp. 14.7.1-14.7.4.
- [34] R. D. Koilpillai and P. P. Vaidyanathan, "Cosine-modulated FIR filter banks satisfying perfect reconstruction," *IEEE Trans. Signal Processing*, vol. 40, no. 4, pp. 770-783, Apr. 1992.
- [35] D. G. Luenberger, *Introduction to Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley, 1973.
- [36] V. P. Sathe and P. P. Vaidyanathan, "Effects of multirate systems on the statistical properties of random signals," *IEEE Trans. Signal Processing*, vol. 41, no. 1, pp. 131-146, Jan. 1993.
- [37] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge University Press, 1985.
- [38] P. J. Davis, *Circulant Matrices*. New York: Wiley, 1979.
- [39] A. N. Akansu and Y. Liu, "On signal decomposition techniques," *Opt. Eng.*, vol. 30, pp. 912-920, July 1991.
- [40] R. C. Agarwal and J. W. Cooley, "New algorithms for digital convolution," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 25, pp. 392-410, Oct. 1977.
- [41] B. Chitprasert and K. R. Rao, "Discrete cosine transform filtering," *Signal Processing*, vol. 19, pp. 233-245, Mar. 1990.
- [42] H. Harada, "On the convolution properties of DCT's and DST's," in *Proc. Int. Symp. Inform. Theory Its Appl.*, Hawaii, Nov. 1990, pp. 591-594.
- [43] K. R. Rao and P. Yip, *Discrete Cosine Transform: Algorithms, Advantages, and Applications*. New York: Academic, 1990.
- [44] S. A. Martucci, "Symmetric convolution and the discrete sine and cosine transforms," preprint.
- [45] T. Chen, I. Djokovic, and P. P. Vaidyanathan, "Results on multi-dimensional nonuniform rational maximally decimated filter banks with orthonormal filters," in *Proc. 26th Annu. Asilomar Conf. Signals, Syst., Comput.*, Oct. 1992.
- [46] J. J. Shynk, "Frequency-domain and multirate adaptive filtering," *IEEE Signal Processing Mag.*, vol. 9, pp. 14-37, Jan. 1992.
- [47] A. Gilloire and M. Vetterli, "Adaptive filtering in subbands with critical sampling: Analysis, experiments, and applications to acoustic echo cancellation," *IEEE Trans. Signal Processing*, vol. 40, pp. 1862-1875, Aug. 1992.
- [48] S. M. Kay and S. L. Marple, "Spectrum analysis: A modern perspective," *Proc. IEEE*, vol. 69, pp. 1380-1419, Nov. 1981.
- [49] J. C. Darragh, "Subband and transform coding of images," doctoral dissertation, Univ. California, Los Angeles, 1989.
- [50] A. Steffen, *Digital Pulse Compression Using Multirate Filter Banks*. Hartung-Gorre Verlag, 1991.



P. P. Vaidyanathan (S'80-M'83-SM'88-F'91) was born in Calcutta, India, on October 16, 1954. He received the B.Sc. (Hons.) degree in physics and the B.Tech. and M.Tech. degrees in radiophysics and electronics, from the University of Calcutta, India, in 1974, 1977, and 1979, respectively, and the Ph.D. degree in electrical and computer engineering from the University of California, Santa Barbara, in 1982.

From 1982 to 1983 he was a postdoctoral fellow at the University of California, Santa Barbara. In 1983 he joined the Electrical Engineering Department of the California Institute of Technology as an Assistant Professor, and since 1988 has been an Associate Professor of Electrical Engineering there. His main research interests are in digital signal processing, multirate systems, wavelet transforms, and adaptive filtering.

Dr. Vaidyanathan served as Vice-Chairman of the Technical Program Committee for the 1983 IEEE International Symposium on Circuits and Systems, and as an Associate Editor for the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS for the period of 1985-1987. He also served as the Technical Program Chairman for the 1992 IEEE International Symposium on Circuits and Systems. He was a recipient of the Award for Excellence in Teaching at the California Institute of Technology for the year of 1983-1984. He also received the NSF's Presidential Young Investigator Award in 1986. In 1989 he received the IEEE ASSP Senior Award for his paper on multirate perfect-reconstruction filter banks. In 1990 he received the S. K. Mitra Memorial Award from the Institute of Electronics and Telecommunications Engineers, India, for his joint paper in the *IETE Journal*.

He was also the coauthor of a paper on linear-phase perfect reconstruction filter banks in the IEEE TRANSACTIONS ON SIGNAL PROCESSING, for which the first author (Truong Nguyen) received the Young Outstanding Author Award in 1993.