# CAUSAL FIR FILTER BANKS WITH ARBITRARY SYSTEM DELAY

*Gerald Schuller* [1]

Bell-Laboratories
600 Mountain Ave.
Murray Hill, NJ07974
email: schuller@lucent.com

*Tanja Karp*

University of Mannheim
B6, 26
68131 Mannheim, Germany
email: karp@rumms.uni-mannheim.de

## ABSTRACT

A design method for causal bi-orthogonal PR FIR $M$-band filter banks is described, which allows an explicit control over system delay, independent of the filter length, with the lowest possible delay equal to the blocking delay of $M - 1$ samples. The design method is very general and can be applied to non-uniform filter banks but also treats uniform modulated filter banks as a special case.

## 1. INTRODUCTION

Filter banks are used in much of the same way as block transforms, for producing a short time frequency domain representation with $M$ frequency bands of a signal (analysis), and reconstructing the signal from this representation (synthesis), but they are a more general approach.

For time-domain data the cascade of the analysis and synthesis filter bank introduces a certain delay, called the system delay. For spatial domain data this is a spatial shift. The system delay is an important property of filter banks. It has not received some attention but recently [1, 2, 3]. Traditionally, filter banks were designed to be orthogonal. In this case the system delay is connected to the filter length: if $L$ is the length of each filter in the analysis and synthesis filter bank, the standard system delay of orthogonal filter banks is $L - 1$ samples. However, many applications require analysis filters with a high stopband attenuation and a small transition bandwidth and thus long analysis filters (if we restrict ourselves to the case of FIR filters) as well as a short system delay. This problem has been partly overcome with the design of bi-orthogonal filter banks [1, 2, 3] where analysis and synthesis filters do not need to be time reversed versions of each other. In this paper we give a formulation that allows for the design of general $M$-band bi-orthogonal filter banks with perfect reconstruction (PR) as well as bi-orthogonal modulated filter banks which are known for their

---

[1] Work was mainly done while with the University of Hanover, Hanover, Germany

low implementation cost, where the system delay can be chosen independently of the filter length.

## 2. DEFINITIONS

For an $M$-band analysis/synthesis filter bank, the input is represented by an $M$-dimensional vector $\mathbf{x}(m)$ composed of the downsampled input components

$$\mathbf{x}(m) = [x(mM+M-1), x(mM+M-2), \ldots, x(mM)]^T$$

Its $z$-transform is the vector $\mathbf{X}(z)$. The polyphase representation for an $M$-band filter bank with input signal $\mathbf{X}(z)$, the subband signal $\mathbf{Y}(z)$, and the reconstructed signal $\hat{\mathbf{X}}(z)$ is

$$\mathbf{Y}(z) = \mathbf{E}(z) \cdot \mathbf{X}(z)$$

for the analysis and

$$\hat{\mathbf{X}}(z) = \mathbf{R}(z) \cdot \mathbf{Y}(z)$$

for the synthesis [4]. $\mathbf{E}(z)$ is the analysis polyphase matrix, $\mathbf{R}(z)$ the synthesis polyphase matrix. Causal filters have no taps at times before zero. This means that $\mathbf{E}(z)$ and $\mathbf{R}(z)$ contain no positive powers of $z$. The filter bank is PR if $\mathbf{R}(z) = z^{-d} \cdot \mathbf{S}^{n_t}(z) \cdot \mathbf{E}^{-1}(z)$ where $\mathbf{S}$ is a Shift Matrix, which circularly shifts the elements of a vector or matrix by one sample [5, 9],

$$\mathbf{S}(z) := \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & \vdots \\ \vdots & \vdots & 0 & \ddots & 0 \\ 0 & & \vdots & & 1 \\ z & 0 & 0 & \cdots & 0 \end{bmatrix}$$

The delay here consists of $d$ blocks of length $M$ minus $n_t$ shifts of single samples. The system delay contains an additional so-called blocking delay of $M - 1$ samples, which results from blocking the input samples into blocks of length $M$ in the polyphase formulation [4]. Hence the

system delay is $n_d = M - 1 + d \cdot M - n_t$. The minimum possible system delay when using causal filters is the blocking delay of $M - 1$ samples, i.e. $n_t$ must be in the range $n_t \leq d \cdot M$. For the filter design it is convenient to limit the range of $n_t$ to $0 \leq n_t \leq M$, without loss of generality.

## 3. THE FACTORIZATION

For the design of filter banks with an arbitrary system delay a formulation based on a factorization of the $M \times M$ polyphase matrices $\mathbf{E}(z)$ and $\mathbf{R}(z)$ into a cascade of a few types of $M \times M$ simple "Filter Matrices" is used. It is a generalization of the formulation for modulated filter banks in [9, 5]. The filter design consists of choosing an appropriate number and type of these filter matrices for the cascade, and then optimizing the coefficients of these matrices for the desired magnitude response.

This decomposition has several advantages for the design. First, it results in bi-orthogonal filter banks where the system delay can be chosen independently of the filter length, contrary to orthogonal filter banks. The perfect reconstruction property and the system delay are structurally guaranteed. Thus, we can perform an unconstrained optimization on the free parameters. Furthermore this cascade can be used to design FIR analysis and synthesis filters ensuring stability. If used for modulated filter banks the factorization also leads to an efficient implementation.

The factorization is based on 2 types of matrices of polynomials of order 1. The first type are "Zero-Delay Matrices", which have the characterizing property that their inverses are causal.

The general Zero-Delay Matrix now has the form

$$\mathbf{L}(z) = \mathbf{I} + \mathbf{A} \cdot z^{-1}$$

where $\mathbf{A}$ is an $M \times M$ matrix with the property

$$\mathbf{A} \cdot \mathbf{A} = \mathbf{0}.$$

Its inverse has the form

$$\mathbf{L}^{-1}(z) = \mathbf{I} - \mathbf{A} \cdot z^{-1}$$

and is causal. Vice versa, if any matrix and its inverse have the form

$$\mathbf{L}(z) = \mathbf{I} + \mathbf{A} \cdot z^{-1}$$
$$\mathbf{L}^{-1}(z) = \mathbf{I} + \mathbf{B} \cdot z^{-1} \tag{1}$$

it follows that

$$\mathbf{B} = -\mathbf{A}$$

and

$$\mathbf{A} \cdot \mathbf{A} = \mathbf{0}$$

The second type of matrix are "Maximum-Delay Matrices" which have an anti-causal inverse, i.e. they contain only non-negative powers of $z$. They have to be multiplied with $z^{-2}$ to make them causal. The general Maximum-Delay Matrix has the form

$$\mathbf{H}(z) = \mathbf{I} \cdot z^{-1} + \mathbf{A}.$$

Its inverse, multiplied with $z^{-2}$, is

$$\mathbf{H}^{-1}(z) \cdot z^{-2} = \mathbf{I} \cdot z^{-1} - \mathbf{A}$$

The same reasoning as in (1) also applies here.

The two matrix types can be used in a product or cascade to construct polyphase matrices with the desired system delay and filter length. To obtain the most general formulation, the cascade also needs a general invertible (real or complex) matrix $\mathbf{T}$, and the Shift Matrix $\mathbf{S}$, which is used for "fine tuning" of the system delay in single steps of the input sampling rate. Hence the polyphase matrices can be written as

$$\mathbf{E}(z) = \mathbf{T} \cdot \prod_{i=1}^{\nu} \mathbf{L}_i(z) \prod_{j=1}^{\mu} \mathbf{H}_j(z) \cdot \mathbf{S}^{n_a}(z)$$

$$\mathbf{R}(z) = \mathbf{S}^{n_s}(z) \cdot \prod_{j=\mu}^{1}(\mathbf{H}_j^{-1}(z) \cdot z^{-2}) \prod_{i=\nu}^{1} \mathbf{L}_i^{-1}(z) \cdot \mathbf{T}^{-1}$$

where each filter matrix $\mathbf{L}_i(z)$ and $\mathbf{H}_j(z)$ has a different matrix $\mathbf{A}$. Because the inverse of the Maximum-Delay Matrices are associated with a multiplication of $z^{-2}$, i.e. a delay of 2 blocks of length $M$, their number $\mu$ determines the system delay, together with the exponents $n_a$ and $n_s$ of the Shift Matrix. Since the system delay additionally contains the blocking delay of $M - 1$ samples it results to

$$n_d = M - 1 + \mu \cdot 2 \cdot M - n_a - n_s$$

The matrices $\mathbf{S}^{n_a}(z)$ and $\mathbf{S}^{n_s}(z)$ are non-causal for $n_a > 0$ or $n_s > 0$. To obtain causal filters their product with $\mathbf{H}_\mu(z)$ and $\mathbf{H}_\mu^{-1}(z) \cdot z^{-2}$ should be causal. This can be obtained by setting the appropriate rows and columns of its $\mathbf{A}$ to zero. These zero-valued rows and columns also reduce the resulting filter length by $\max(n_a, n_s)$. The matrix $\mathbf{T}$ alone leads to filters of length $M$. Each Filter Matrix increases the filter length by $M$, and the Shift-Matrices reduce the filter length by $\max(n_a, n_s)$. This means the resulting filter length is mainly determined by the total number of Filter Matrices $\mu + \nu$ to

$$L = (\mu + \nu) \cdot M + M - \max(n_a, n_s)$$

Since more entries of the filter matrices can be zero, this is the maximum filter length. For cosine-modulated filter banks this kind of decomposition leads to sparse matrices and hence an efficient implementation, similar to the decomposition described in [5, 1]. In this case $\mathbf{T}$ is a Discrete

Cosine Transform (e.g. a DCT-4), multiplied from the right with a diagonal matrix, and the $\mathbf{A}$ matrices are zero except for one half of the anti-diagonal. Observe that this cascade is also suitable for an implementation with low computational accuracy, since the inverse filter matrices result from sign flipping, so that PR is maintained even after coefficient quantization.

In general, the matrices $\mathbf{A}$ have rank $\leq M/2$ (since $\mathbf{A} \cdot \mathbf{A} = 0$) so they can be written as $\mathbf{A} = \mathbf{v} \cdot \mathbf{w}$ where $\mathbf{v}$ is a $M \times M/2$ and $\mathbf{w}$ is a $M/2 \times M$ matrix. If $\mathbf{v}$ and $\mathbf{w}$ are orthogonal, i.e. if $\mathbf{w} \cdot \mathbf{v} = 0$ then

$$\mathbf{A} \cdot \mathbf{A} = \mathbf{v} \cdot (\mathbf{w} \cdot \mathbf{v}) \cdot \mathbf{w} = 0$$

which is the desired property.

To determine the degrees of freedom of one matrix $\mathbf{A}$, it can be observed that e.g. $\mathbf{w}$ can be normalized (i.e. it consists of unit norm row vectors of length $M$), since any factor can be made part of $\mathbf{v}$. The rows of $\mathbf{w}$ can also be ordered, because the effect of any reordering of its rows can be obtained by reordering the columns of $\mathbf{v}$. This means that all degrees of freedom are in $\mathbf{v}$, which contains $M \cdot M/2$ degrees of freedom.

For example for the 2-band case, $M = 2$, the orthogonal vectors $\mathbf{v}$ and $\mathbf{w}$ can be written as

$$\mathbf{v} = a \cdot [\sin(\alpha), \cos(\alpha)]^T$$

$$\mathbf{w} = [-\cos(\alpha), \sin(\alpha)]$$

where $a$ and $\alpha$ are then the unknowns which determine $\mathbf{A}$. The coefficients of the cascade are obtained by optimizing them for a desired magnitude response. It can be observed, that neighboring filter matrices should have different $\mathbf{A}$ matrices, so that they increase the filter length.

The order of the Maximum-Delay and Zero-Delay Matrices is not important, as can be seen in the following proof of effective completeness. This proof shows, that all FIR filter banks with perfect reconstruction can be approximated arbitrary close by this cascade (hence "effective" completeness). It is similar to the one described in [5] for the cosine-modulated case, with some modifications for the general case. Assume an FIR PR filter bank is given, with filter length $L$ and system delay $n_d$. First define $\mathbf{F_E}(z)$ and $\mathbf{F_R}(z)$ as the part of the cascade without the shift matrix, obtained by

$$\mathbf{F_E}(z) = \mathbf{E}(z) \cdot \mathbf{S}^{-n_a}(z)$$

$$\mathbf{F_R}(z) = \mathbf{S}^{-n_s}(z) \cdot \mathbf{R}(z)$$

where $n_a$ and $n_s$ are chosen such that

$$n_d = M - 1 + d \cdot M - n_a - n_s$$

for some even integer $d$, and such that $0 \leq n_a = n_s \leq M$ (for even delays). Then write, with $K = L/M$ (with appropriate rounding for $L/M$ if necessary)

$$\mathbf{F_E}(z) = \sum_{m=0}^{K-1} \mathbf{f_E}(m) \cdot z^{-m}$$

$$\mathbf{F_R}(z) = \sum_{m=0}^{K-1} \mathbf{f_R}(m) \cdot z^{-m}$$

PR and the system delay lead to $\mathbf{F_R}(z) \cdot \mathbf{F_E}(z) = z^{-d} \cdot \mathbf{I}$. By analyzing the sums it can be seen that for $d < 2L/M - 3$ this means that

$$(\mathbf{f_R}(K-2)+\mathbf{f_R}(K-1)\cdot z^{-1})\cdot(\mathbf{f_E}(K-2)+\mathbf{f_E}(K-1)\cdot z^{-1}) =$$
$$= \mathbf{f_R}(K-2) \cdot \mathbf{f_E}(K-2)$$

and a Zero-Delay matrix can be extracted by setting

$$\mathbf{L}_i(z) = \mathbf{I} + \mathbf{f_E}(K-1) \cdot (\mathbf{f_E}(K-2))^{-1} \cdot z^{-1}$$

$$\mathbf{L}_i^{-1}(z) = \mathbf{I} + (\mathbf{f_R}(K-2))^{-1} \cdot \mathbf{f_R}(K-1) \cdot z^{-1}$$

where $i$ is initially $i = 1$. Since they have the desired form of (1), it means that

$$\mathbf{A} = \mathbf{f_E}(K-1) \cdot (\mathbf{f_E}(K-2))^{-1} =$$
$$= -\mathbf{f_R}(K-2))^{-1} \cdot \mathbf{f_R}(K-1)$$

and $\mathbf{A} \cdot \mathbf{A} = 0$, the required property. This extraction is only possible, if $\mathbf{f_E}(K-2)$ and $\mathbf{f_R}(K-2)$ are invertible. But if they are not invertible, some arbitrary small $\epsilon$ can be added to them to make them invertible. This way an arbitrary close approximation is possible. In numerical simulations this $\epsilon$ may lead to problems because of near singular matrices, but it shows that a close approximation is possible. After extracting the Zero-Delay Matrix the matrices $\mathbf{F_E}(z)$ and $\mathbf{F_R}(z)$ are replaced by $\mathbf{L}_i^{-1}(z)\mathbf{F_E}(z)$ and $\mathbf{F_R}(z)\mathbf{L}_i(z)$, which have filters with lengths reduced by $M$. For the next extraction $i$ is increased by one to $i = 2$. This can be repeated until $d \geq 2L - 3$.

The Maximum-Delay Matrices can be extracted in a similar manner, as long as $d > 1$, by setting

$$\mathbf{H}_j(z) = \mathbf{I} \cdot z^{-1} + \mathbf{f_E}(0) \cdot (\mathbf{f_E}(1))^{-1}$$

$$\mathbf{H}_j^{-1}(z) \cdot z^{-2} = \mathbf{I} \cdot z^{-1} + (\mathbf{f_R}(1))^{-1} \cdot \mathbf{f_R}(0)$$

The remaining matrix is the real or complex matrix $\mathbf{T}$ or $\mathbf{T}^{-1}$ resp. But here it appears on the right side of $\mathbf{E}(z)$ and the left side of $\mathbf{R}(z)$ as in the case of GenLOT [10] which, however, only treats the orthogonal case. Especially for modulated filter banks it is advantageous to have it in the opposite side, since it leads to sparser matrices and hence a more efficient implementation. This can simply be achieved by replacing the extracted matrices $\mathbf{L}_i(z)$ by $\mathbf{T}^{-1}\mathbf{L}_i(z)\mathbf{T}$ and $\mathbf{H}_i(z)$ by $\mathbf{T}^{-1}\mathbf{H}_i(z)\mathbf{T}$. For modulated filter banks this

approach can also be used to find the form of the sparse filter matrices, i.e. the positions where the non-zero elements are located, e.g. by using an example filter bank. Since this extraction works for any order of the Maximum-Delay and Zero-Delay Matrices, this also shows, that their ordering is not important in principle. But it may still be important for a numerical implementation.

The general approach results in filter matrices which are not as sparse as for the modulated case, so that an implementation will not be as efficient as for the modulated case. But the increased number of coefficients has the advantage, that it also means an increased number in the degrees of freedom for the design process. This is especially important for the case of filter banks with only a few bands, e.g. for the 2-band case. Other applications are e.g. non-uniform filter banks. The frequency responses are obtained by

$$\mathbf{E}(z) \cdot [1, \ldots, z^{-(M-1)}]^T$$

for the analysis, and

$$[z^{-(M-1)}, \ldots, 1] \cdot \mathbf{R}(z)$$

for the synthesis [4]. For the 2-band case it can be useful to have a zero at $\omega = 0$ for the higher band. This can be obtained by setting $z = 1$ in above formulas and setting the response of the higher band to zero. This can then be used as a constraint for the optimization of the frequency response.

For the 2-band case a similarity to the methods in [6, 7, 8] can be observed. E.g. the lifting-scheme described in [8] uses only 2 "steps" or matrices in the cascade, but of arbitrary order, and [7] describes a factorization but no design method. In contrast, the presented scheme uses arbitrary many matrices with fixed order, which are directly connected to the system delay.

## 4. DESIGN EXAMPLE

Figure 1 shows an example for the 2-band case. The matrix coefficients were obtained with the optimization method described in [1]. It shows a comparison of a 2-band filter bank with a standard system delay, and a filter bank with the same system delay but longer filters (hence a low delay filter bank). Both where designed with the described algorithm. The low delay filter bank also has a magnitude response similar to that of the QMF filter bank used in the G.722 speech coder, but a lower system delay (9 samples vs. 23 samples), so that its use could reduce the coding delay of that coder.

## 5. REFERENCES

[1] G.D.T. Schuller and M. J. T. Smith, "New Framework for Modulated Perfect Reconstruction Filter Banks", IEEE Trans. on Sig. Proc., Aug. 1996.
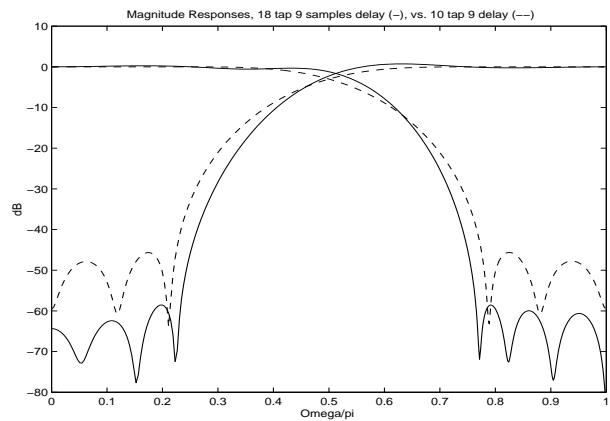
Figure 1: Analysis magnitude responses for a 2-band low delay filter bank with 18 taps and 9 samples system delay (solid line), compared with a filter bank with a standard delay, with 10 taps and the same system delay of 9 samples (dashed line).

[2] P.N. Heller, T. Karp, and T.Q. Nguyen, "A General Formulation of Modulated Filter Banks", subm. to IEEE Trans. on Sig. Proc., 1996

[3] S.-M. Phoong, C.W. Kim, P.P. Vaidyanathan, and R. Ansari, "A New Class of Two-Channel Biorthogonal Filter Banks and Wavelet Bases", IEEE Trans. on Sig. Proc., March 1995.

[4] P.P. Vaidyanathan, *Multirate Systems and Filter Banks*, Prentice Hall, 1993.

[5] G. Schuller, "A New Factorization and Structure for Cosine Modulated Filter Banks with Variable System Delay", Asilomar Conf. on Sign., Syst., and Computers, Pacific Grove, CA, Nov. 3-6, 1996

[6] A.A.C. Kalker, I.A. Shah, "Ladder Structures for Multidimensional Linear Phase Perfect Reconstruction Filter Banks and Wavelets", Visual Communications and Image Processing '92, pp. 12-20.

[7] I. Daubechies, W. Sweldens, "Factoring Wavelet Transforms into Lifting Steps", Preprint, Bell Laboratories, Lucent Technologies, 1996.

[8] W. Sweldens, "The Lifting Scheme: A Custom-Design Construction of Biorthogonal Wavelets", Applied and Computational Harmonic Analysis 3, 186-200, 1996.

[9] G. Schuller, "Time-Varying Filter Banks with Variable System Delay" ICASSP 97, Munich, Germany, Apr. 1997.

[10] R. L. de Queiroz, T. Q. Nguyen, and K. R. Rao, "The GenLOT: Generalized Linear-Phase Lapped Transform" IEEE Trans. on Sig. Proc., March 1996.