

# Optimized Signal Expansions for Sparse Representation

Sven Ole Aase, John Håkon Husøy, Karl Skretting, and Kjersti Engan

**Abstract**—Traditional signal decompositions such as transforms, filterbanks, and wavelets generate signal expansions using the *analysis–synthesis* setting: The expansion coefficients are found by taking the inner product of the signal with the corresponding analysis vector. In this paper, we try to free ourselves from the analysis–synthesis paradigm by concentrating on the *synthesis* or *reconstruction* part of the signal expansion. Ignoring the analysis issue completely, we construct sets of synthesis vectors, which are denoted *waveform dictionaries*, for efficient signal representation. Within this framework, we present an algorithm for designing waveform dictionaries that allow *sparse* representations: The objective is to approximate a training signal using a small number of dictionary vectors. Our algorithm optimizes the dictionary vectors with respect to the average nonlinear approximation error, i.e., the error resulting when keeping a fixed number  $n$  of expansion coefficients but not necessarily the *first*  $n$  coefficients. Using signals from a Gaussian, autoregressive process with correlation factor 0.95, it is demonstrated that for established signal expansions like the Karhunen–Loève transform, the lapped orthogonal transform, and the biorthogonal 7/9 wavelet, it is possible to improve the approximation capabilities by up to 30% by fine tuning of the expansion vectors.

## I. INTRODUCTION

A *signal expansion* is simply a weighted sum of vectors  $\mathbf{f}_i$ . This weighted sum may be identical to, or an approximation to, a given signal vector  $\mathbf{x}$ . If the expansion is identical to  $\mathbf{x}$ , we can write

$$\mathbf{x} = \mathcal{F}\mathbf{w} = \sum_i w_i \mathbf{f}_i \quad (1)$$

where  $\mathcal{F}$  is a (possibly infinite) matrix with  $\{\mathbf{f}_i\}$  as columns, and  $\mathbf{w}$  is the vector of expansion coefficients. Equation (1) can be interpreted as a *synthesis* formula in the sense that  $\mathbf{x}$  is synthesized, or build up, from a library of expansion vectors using appropriately selected values for the expansion coefficients. For this reason,  $\mathcal{F}$  is sometimes referred to as a *waveform dictionary*. If the matrix  $\mathcal{F}$  is invertible, a unique set of coefficients for the exact representation of any signal vector  $\mathbf{x}$  can be obtained as<sup>1</sup>

$$\mathbf{w} = \mathcal{F}^{-1}\mathbf{x} \quad (2)$$

Manuscript received September 24, 1999; revised January 9, 2001. The associate editor coordinating the review of this paper and approving it for publication was Prof. Lang Tong.

The authors are with Høgskolen i Stavanger, Department of Electrical and Computer Engineering, Stavanger, Norway (e-mail: Sven.O.Aase@tn.his.no).  
 Publisher Item Identifier S 1053-587X(01)03338-4.

<sup>1</sup>For notational convenience, we denote the *forward* matrix by  $\mathcal{F}^{-1}$  and the *inverse* or *reconstruction* matrix by  $\mathcal{F}$ .

and this is commonly referred to as the *analysis* equation in an analysis–synthesis setting. Note that in the interest of maximum generality, we have not specified the dimensions of the matrices and vectors involved in (1) and (2). Depending on the dimensions, which may extend to infinity, as well as the structure of the  $\mathcal{F}$  matrix, the analysis–synthesis equations given above cover many important cases including transforms, filterbanks, wavelets, and wavelet packets.

The main objective for using the analysis–synthesis framework in signal processing applications is to construct  $\mathcal{F}$  such that the vector of coefficients  $\mathbf{w}$  is more attractive to work with than  $\mathbf{x}$ . For signal representation purposes, a crucial point is that  $\mathbf{w}$  should be *sparse*. The sparseness constraint refers to the requirement that  $\mathbf{w}$  must have as few nonzero elements as possible [1]. In a data compression setup, the sparseness constraint facilitates bit-efficient representation of the original vector  $\mathbf{x}$  since only the nonzero elements of  $\mathbf{w}$  have to be quantized and stored or transmitted.

In the present work, we aim at freeing ourselves from the traditional analysis–synthesis paradigm in that we concentrate on the *synthesis* or *reconstruction* part of the signal expansion. That is, given the coefficients  $\{w_i\}$ , the reconstructed signal vector  $\tilde{\mathbf{x}}$  is given by

$$\tilde{\mathbf{x}} = \mathcal{F}\mathbf{w} = \sum_i w_i \mathbf{f}_i. \quad (3)$$

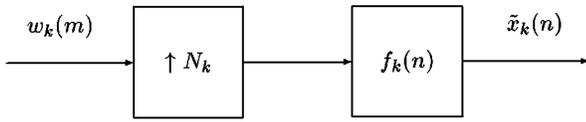
The consequences of ignoring the analysis part of the expansion can be summarized as follows.

- 1) There are no restrictions on the choice of waveform dictionary  $\mathcal{F}$ . The invertibility of  $\mathcal{F}$  is no longer an issue, and  $\mathcal{F}$  may be *overcomplete*. This means that the number of columns, or dictionary vectors, may be larger than the dimension of  $\mathbf{x}$ . The motivation for allowing overcomplete dictionaries is simply to increase the choice of possible expansion vectors in the hope that this allows for better approximation capability while preserving the sparseness criterion.
- 2) Since we no longer assume that  $\mathcal{F}$  is neither orthogonal nor invertible, we obviously have to find some other way of determining the expansion coefficients than what is done in traditional signal decompositions. We have employed the *fast orthogonal matching pursuit* (FOMP) algorithm of [2].

The main issue dealt with in this paper is as follows: Given a training signal denoted  $\mathbf{x}$ , or  $x(n)$ , *how do we find the waveform dictionary such that the best possible signal approximation results when keeping a small collection of coefficients?* In a practical setting, the training signal would contain different





Fig. 2. Synthesis from channel number  $k$ .

should allow the MSE to grow for several iterations without terminating the training. This can be seen from a training example in Section III-A.

### B. Computing the Optimal Dictionary Update

We now explain how to perform the improvement in  $\mathcal{F}$  constituting step 3 of the design algorithm. The objective is to find the  $\mathcal{F}$  that minimizes the approximation residual while keeping the already-computed coefficients in the expansion resulting from step 2 (or 4) of the previous iteration of the algorithm. In the derivations to follow, the optimal  $\mathcal{F}$  is found by simply differentiating the quadratic approximation error with respect to all elements of the dictionary vectors. For all cases considered, it is shown that this leads to a linear system of equations solvable for  $\mathcal{F}$ .

We start by solving the problem in its full generality using filter bank terminology. As explained in the introduction of Section II, (6) and (7) can be viewed as a multichannel synthesis system, where  $\mathcal{F}_k \mathbf{w}_k$  denote the contribution from channel number  $k$ . This is illustrated in Fig. 2, where  $w_k(m)$  are the elements of the vector  $\mathbf{w}_k$  in (6), which is now interpreted as a time series.

The contribution from channel number  $k$  to the synthesized signal can be written [13]

$$\tilde{x}_k(n) = \sum_{p=0}^{\lfloor (L_k-1)/N_k \rfloor} w_k \left( \left\lfloor \frac{n}{N_k} \right\rfloor - p \right) f_k(n \bmod N_k + pN_k) \quad (11)$$

where  $\lfloor y \rfloor$  denotes the largest integer value less than, or equal to,  $y$ . The reconstructed signal is

$$\tilde{x}(n) = \sum_{k=0}^{K-1} \tilde{x}_k(n). \quad (12)$$

We define the object function  $J$  as

$$J = \sum_{n=-\infty}^{\infty} (x(n) - \tilde{x}(n))^2 \quad (13)$$

and by setting its gradient to zero a straightforward derivation leads to

$$\begin{aligned} & \sum_{n=-\infty}^{\infty} \sum_{k=0}^{K-1} \sum_{p=0}^{\lfloor (L_k-1)/N_k \rfloor} \sum_{t=0}^{K-1} \sum_{s=0}^{\lfloor (L_t-1)/N_t \rfloor} f_t(n \bmod N_t + sN_t) \\ & \cdot w_k \left( \left\lfloor \frac{n}{N_k} \right\rfloor - p \right) w_t \left( \left\lfloor \frac{n}{N_t} \right\rfloor - s \right) \\ & \cdot \frac{\partial f_k(n \bmod N_k + pN_k)}{\partial f_r(q)} \\ & = \sum_{n=-\infty}^{\infty} \sum_{k=0}^{K-1} \sum_{p=0}^{\lfloor (L_k-1)/N_k \rfloor} x(n) w_k \left( \left\lfloor \frac{n}{N_k} \right\rfloor - p \right) \\ & \cdot \frac{\partial f_k(n \bmod N_k + pN_k)}{\partial f_r(q)}. \end{aligned} \quad (14)$$

Here, we have differentiated with respect to  $f_r(q)$ ,  $r = 0, \dots, K-1$ , and  $q = 0, \dots, L_r-1$ .

The derivative term in (14) vanishes for all values of  $n, k$ , and  $p$ , except when

- 1)  $k = r$ ;
- 2)  $n \bmod N_r + pN_r = q$

in which case, it equals unity. Using the division algorithm due to Euclid, we may decompose  $q$  as

$$q = \left\lfloor \frac{q}{N_r} \right\rfloor N_r + q \bmod N_r. \quad (15)$$

Comparing (15) with the expression in condition number 2 above, the uniqueness of the decomposition enforces  $p = \lfloor q/N_r \rfloor$  and  $n = q \bmod N_r + mN_r$ , where  $m$  is a new summation variable. Equation (14) now becomes

$$\begin{aligned} & \sum_{t=0}^{K-1} \sum_{s=0}^{\lfloor \frac{L_t-1}{N_t} \rfloor} \sum_{m=-\infty}^{\infty} f_t((q \bmod N_r + mN_r) \bmod N_t + sN_t) \\ & \cdot w_r \left( m - \left\lfloor \frac{q}{N_r} \right\rfloor \right) w_t \left( \left\lfloor \frac{q \bmod N_r + mN_r}{N_t} \right\rfloor - s \right) \\ & = \sum_{m=-\infty}^{\infty} x(q \bmod N_r + mN_r) w_r \left( m - \left\lfloor \frac{q}{N_r} \right\rfloor \right). \end{aligned} \quad (16)$$

*Left Side:* We now focus on the left-hand side of (16), where the goal is to remove the infinite summation over the  $f_t(\cdot)$  factor with the objective of establishing a finite, linear set of equations in the  $f_t(\cdot)$ s. We change the summation variable as

$$m = Nu + v, \quad \text{where } -\infty < u < \infty, 0 \leq v < N \quad (17)$$

where  $N$  is the least common multiple of the upsampling factors:

$$N = \text{LCM}(N_0, \dots, N_{K-1}). \quad (18)$$

Using (17), we can write

$$\begin{aligned} & \sum_{v=0}^{N-1} \sum_{u=-\infty}^{\infty} f_t((q \bmod N_r + uN_rN + vN_r) \bmod N_t + sN_t) \\ & \cdot w_r \left( uN + v - \left\lfloor \frac{q}{N_r} \right\rfloor \right) \\ & \cdot w_t \left( \left\lfloor \frac{q \bmod N_r + uN_rN + vN_r}{N_t} \right\rfloor - s \right) \\ & = \sum_{v=0}^{N-1} f_t((q \bmod N_r + vN_r) \bmod N_t + sN_t) \\ & \cdot \sum_{u=-\infty}^{\infty} w_r \left( uN + v - \left\lfloor \frac{q}{N_r} \right\rfloor \right) \\ & \cdot w_t \left( \left\lfloor \frac{q \bmod N_r + vN_r}{N_t} \right\rfloor + uN_r \frac{N}{N_t} - s \right). \end{aligned} \quad (19)$$

*Right Side:* Using the fact that

$$\sum_{m=-\infty}^{\infty} g(m+i)h(pm) = \sum_{m=-\infty}^{\infty} g(m)h(p(m-i)) \quad (20)$$

the expression on the right-hand side in (16) can be simplified as (21)

$$\begin{aligned} & \sum_{m=-\infty}^{\infty} x\left(q \bmod N_r + \left\lfloor \frac{q}{N_r} \right\rfloor N_r + mN_r\right) w_r(m) \\ &= \sum_{m=-\infty}^{\infty} x(q + mN_r) w_r(m). \end{aligned} \quad (21)$$

*Putting it All Together:* Collecting the results from (19) and (21), the optimal values  $\{f_r(q)\}$  must satisfy

$$\begin{aligned} & \sum_{t=0}^{K-1} \sum_{s=0}^{\lfloor \frac{L_t-1}{N_t} \rfloor} \sum_{v=0}^{N-1} f_t((q \bmod N_r + vN_r) \bmod N_t + sN_t) \\ & \cdot \sum_{u=-\infty}^{\infty} w_r\left(uN + v - \left\lfloor \frac{q}{N_r} \right\rfloor\right) \\ & \cdot w_t\left(\left\lfloor \frac{q \bmod N_r + vN_r}{N_t} \right\rfloor + uN_r \frac{N}{N_t} - s\right) \\ &= \sum_{m=-\infty}^{\infty} x(q + mN_r) w_r(m) \end{aligned} \quad (22)$$

for all values  $r = 0, \dots, K-1$ , and  $q = 0, \dots, L_r - 1$ .

A matrix formulation equivalent to (22) is obtained by setting

$$\begin{aligned} a(r, q, t, s, v) &= \sum_{u=-\infty}^{\infty} w_r\left(uN + v - \left\lfloor \frac{q}{N_r} \right\rfloor\right) \\ & \cdot w_t\left(\left\lfloor \frac{q \bmod N_r + vN_r}{N_t} \right\rfloor + uN_r \frac{N}{N_t} - s\right) \end{aligned} \quad (23)$$

$$b_{rq} = \sum_{m=-\infty}^{\infty} x(q + mN_r) w_r(m) \quad (24)$$

$$g(r, q, t, s, v) = (q \bmod N_r + vN_r) \bmod N_t + sN_t \quad (25)$$

$$A_{ti}^{rq} = \sum_{\{(s,v) | g(r,q,t,s,v)=i\}} a(r, q, t, s, v) \quad (26)$$

$$\begin{aligned} \mathbf{f} &= (f_0(0), f_0(1), \dots \\ & \dots, f_{K-1}(L_{K-1} - 2), f_{K-1}(L_{K-1} - 1))^T \end{aligned} \quad (27)$$

giving

$$\begin{aligned} & \begin{bmatrix} A_{00}^{00} & A_{01}^{00} & \dots & A_{K-1, L_{K-1}-1}^{00} \\ A_{00}^{01} & A_{01}^{01} & \dots & A_{K-1, L_{K-1}-1}^{01} \\ \vdots & \vdots & \ddots & \vdots \\ A_{00}^{K-1, L_{K-1}-1} & A_{01}^{K-1, L_{K-1}-1} & \dots & A_{K-1, L_{K-1}-1}^{K-1, L_{K-1}-1} \end{bmatrix} \\ & \cdot \mathbf{f} = \begin{bmatrix} b_{00} \\ b_{01} \\ \vdots \\ b_{K-1, L_{K-1}-1} \end{bmatrix}. \end{aligned} \quad (28)$$

From this linear set of  $\sum_{r=0}^{K-1} L_r$  equations, we can compute the vector  $\mathbf{f}$  constituting the best possible dictionary  $\mathcal{F}$ , in the mean square error sense, when the coefficients for each

vector approximation are fixed. Since the object function (13) is bounded downwards, we know that the problem has a minimum solution. Since (22) only has one solution when the linear system has full rank, this must be the unique, global solution.

Notice again the generality of (22) and (28). A vast set of different configurations may be optimized using the formulation above, including cases where the dictionary is under or over-complete.

In the following subsections, we study three very important special cases and their generalizations:

- 1) Block transforms/frames;
- 2) uniform filter banks/LOTs
- 3) wavelet-like expansions.

*1) Special Case Number 1: Block Transforms/Frames:* If we set  $N_0 = N_1 = \dots = N_{K-1} = N$  and  $L_0 = L_1 = \dots = L_{K-1} = N$ , we obtain a block-oriented signal expansion. If  $K = N$ , we have a transform, and if  $K > N$ , we have a block-oriented frame.

Equation (22) can now be written as

$$\sum_{t=0}^{K-1} f_t(q) \sum_{m=-\infty}^{\infty} w_r(m) w_t(m) = \sum_{m=-\infty}^{\infty} x(q + mN) w_r(m) \quad (29)$$

where  $r = 0, \dots, K-1$  and  $q = 0, \dots, N-1$ .

Defining the  $N \times K$  matrix  $\mathbf{F}$  as

$$\mathbf{F} = \begin{bmatrix} f_0(0) & f_1(0) & \dots & f_{K-1}(0) \\ f_0(1) & f_1(1) & \dots & f_{K-1}(1) \\ \vdots & \vdots & \ddots & \vdots \\ f_0(N-1) & f_1(N-1) & \dots & f_{K-1}(N-1) \end{bmatrix} \quad (30)$$

and the matrices  $\mathbf{A}$  and  $\mathbf{B}$  by their elements as

$$A_{rt} = \sum_{m=-\infty}^{\infty} w_r(m) w_t(m) \quad (31)$$

$$B_{qr} = \sum_{m=-\infty}^{\infty} x(q + mN) w_r(m) \quad (32)$$

(29) can be written as

$$\mathbf{FA} = \mathbf{B} \quad (33)$$

with the solution

$$\mathbf{F} = \mathbf{BA}^{-1}. \quad (34)$$

This is equivalent to the result derived in [3] and [11].

*2) Special Case Number 2: Uniform Filter Banks/LOTs:* If we set  $N_0 = N_1 = \dots = N_{K-1} = N$  and  $L_0 = L_1 = \dots = L_{K-1} = PN$ , where  $P$  is a positive integer, we obtain a signal expansion similar to a uniform,  $K$ -channel, FIR filterbank where the channels have the same filter length. This is sometimes referred to as a LOT [14], [15].

Using the division algorithm due to Euclid, we can decompose  $q$  as

$$q = q_1 + q_2N \quad (35)$$

where  $q_1 = q \bmod N$ , and  $q_2 = \lfloor q/N \rfloor$ . Equation (22) then becomes

$$\begin{aligned} & \sum_{t=0}^{K-1} \sum_{s=0}^{P-1} f_t(q_1 + sN) \underbrace{\sum_{m=-\infty}^{\infty} w_r(m - q_2) w_t(m - s)}_{a(t,s,q_2,r)} \\ &= \underbrace{\sum_{m=-\infty}^{\infty} x(q_1 + q_2N + mN) w_r(m)}_{b(q_1,q_2,r)} \end{aligned} \quad (36)$$

where  $r = 0, \dots, K-1$ ,  $q_1 = 0, \dots, N-1$ ,  $q_2 = 0, \dots, P-1$ .

Defining the  $N \times PK$  matrix  $\mathbf{F}^*$  as (37), shown at the bottom of the page, and the matrices  $\mathbf{A}$  and  $\mathbf{B}$  by their elements as

$$\begin{aligned} A_{kl} &= a \left( \left\lfloor \frac{k}{P} \right\rfloor, k \bmod P, l \bmod P, \left\lfloor \frac{l}{P} \right\rfloor \right) \\ & k, l = 0, \dots, PK - 1 \\ B_{kl} &= b \left( k, l \bmod P, \left\lfloor \frac{l}{P} \right\rfloor \right) \\ & k = 0, \dots, N - 1, \quad l = 0, \dots, PK - 1 \end{aligned} \quad (38)$$

Equation (36) can be written as

$$\mathbf{F}^* \mathbf{A} = \mathbf{B} \quad (39)$$

with the solution

$$\mathbf{F}^* = \mathbf{B} \mathbf{A}^{-1}. \quad (40)$$

Note that the  $\mathbf{F}^*$  matrix is simply introduced in order to provide a compact formulation for the solution of the optimal dictionary in the uniform FIR filterbank case and does not relate directly to the internal structure of the dictionary matrix itself. We would also like to emphasize the possibility of optimizing configurations where  $K \neq N$ .

3) *Special Case Number 3: Wavelets*: A discrete-time wavelet decomposition is constructed by repeatedly using two filters (lowpass and highpass) on the output of the lowpass filter in each stage. In this manner, a dyadic frequency partitioning results [13].

The general framework outlined in Section II-B can easily be adapted to mimic a general,  $M$ -stage, wavelet-type expansion by setting  $K = M + 1$  and the upsampling parameters as follows:

$$\begin{aligned} N_0 &= 2^M \text{ ``Lowpass channel''} \\ N_1 &= 2^M \text{ ``1. bandpass channel''} \\ N_2 &= 2^{M-1} \text{ ``2. bandpass channel''} \\ N_3 &= 2^{M-2} \text{ ``3. bandpass channel''} \\ &\vdots \\ N_M &= 2 \text{ ``Highpass channel''}. \end{aligned} \quad (41)$$

The optimal filters are found by inserting the chosen upsampling parameters and filter lengths into (28) and setting  $N = 2^M$ . A varied set of under or overcomplete wavelet-like signal expansions can be generated in the same manner. Note that the resulting expansions are not wavelets but expansions mimicking the *support structure* of wavelets.

### C. Finding a Signal Approximation

As explained in the introduction, by giving up the invertibility of the dictionary matrix  $\mathcal{F}$ , we have to find some other way of computing the expansion coefficients of (3). The goal is to use few coefficients while constructing a good approximation to the training signal  $x(n)$  or  $\mathbf{x}$ .

Ideally, when approximating the signal  $x(n)$ , the available coefficients should be allocated over the whole signal. To illustrate this idea, let us focus on the simple case where the dictionary is a block-oriented signal expansion, where the  $N \times K$  matrix  $\mathbf{F}$  is defined as in (30) in Section II-B1. Furthermore, we assume the training signal to be finite, with length  $L_x = MN$ , where  $M$  is a positive integer.

The signal expansion can be written using  $\mathbf{F}$  repeated  $M$  times, as illustrated in

$$\begin{bmatrix} \mathbf{F} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{F} & & \\ \vdots & & \ddots & \\ \mathbf{0} & & & \mathbf{F} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{w}^{(0)} \\ \mathbf{w}^{(1)} \\ \vdots \\ \mathbf{w}^{(M-1)} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{x}}^{(0)} \\ \tilde{\mathbf{x}}^{(1)} \\ \vdots \\ \tilde{\mathbf{x}}^{(M-1)} \end{bmatrix} \quad (42)$$

where  $\mathbf{x}^{(i)} \in R^N$  and  $\mathbf{w}^{(i)} \in R^K$  denote the  $i$ th block of the signal and coefficient vector, respectively.

If the available coefficients are allocated freely among the  $M$   $\mathbf{F}$ -blocks, signal blocks with high energy content can be approximated using more coefficients than blocks having less energy. This corresponds to the practical situation occurring in straight-forward transform coding when the vectors to be retained in the expansion are found by keeping those coefficients above a given threshold. Clearly, if  $M$  is small, this scheme would offer little flexibility for approximating signal regions having varying characteristics.

Given a sparseness factor  $S$  defined as in (9) and with practical considerations in mind, the global procedure of selecting  $L_x S$  out of  $L_x K/N$  vectors is prohibitive when  $L_x$  is large, as will be the case in our experiments. It follows that the procedure of selecting vectors (and correspondingly, expansion coefficients) must be performed in a segment-by-segment manner.

We partition the training signal  $x(n)$  into equally sized segments, and, depending on whether the dictionary is block-based or not, the procedure is performed as follows:

*Block-Based Dictionary*: Using a segment size equal to  $MN$ , each segment of  $x(n)$  will be *independently* approximated, as

$$\mathbf{F}^* = \begin{bmatrix} f_0(0) & f_0(N) & \cdots & f_0((P-1)N) & \cdots & f_{K-1}(0) & \cdots & f_{K-1}((P-1)N) \\ f_0(1) & f_0(N+1) & \cdots & f_0((P-1)N+1) & \cdots & f_{K-1}(1) & \cdots & f_{K-1}((P-1)N+1) \\ & & \vdots & & & \vdots & & \\ f_0(N-1) & f_0(2N-1) & \cdots & f_0(PN-1) & \cdots & f_{K-1}(N-1) & \cdots & f_{K-1}(PN-1) \end{bmatrix} \quad (37)$$

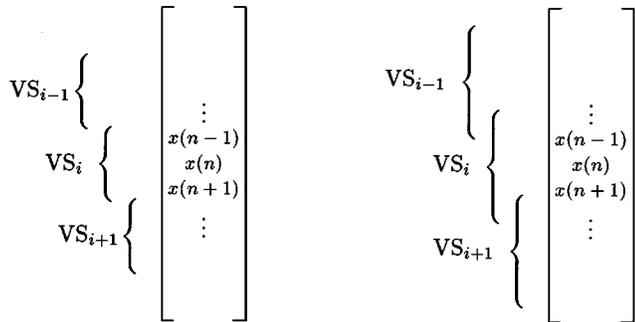


Fig. 3. Vector selection (VS) is performed segment-by-segment according to dictionary type: Block-based (left) and overlap (right).

shown in (42), using  $MNS$  vectors. The procedure is shown on the left side of Fig. 3.

*Dictionary with Overlapping Vectors:* If the dictionary has overlapping vectors (for example, the uniform filterbank case in Section II-B2 or the wavelet case in Section II-B3), a perfect separation of the segment-based expansions is not possible. Due to overlapping vectors, adjacent expansions will overlap.

A practical way of overcoming the overlap problem is as follows: The training signal  $x(n)$  is approximated segment-by-segment, in a top-down manner, as shown on the right side of Fig. 3. When the expansion for the current segment has been found, the associated reconstruction is subtracted from  $x(n)$  before starting the procedure on the next segment.

1) *Vector Selection Algorithms:* For each segment of  $x(n)$ , the task is to find the best possible approximation using a predetermined number of dictionary vectors. If we wanted to find the *optimal* approximation to a given segment by selecting  $a$  out of  $b$  dictionary vectors, it would be necessary to investigate

$$\binom{b}{a} \quad (43)$$

different choices of vectors. Thus, finding the optimal vectors to use in an approximation is an NP-hard problem and requires extensive calculation [16]. It follows that a suboptimal algorithm is preferable in order to limit the computational complexity. There exist several different vector selection schemes for solving this problem [1], [17]–[19]. In this paper, the *fast orthogonal matching pursuit* (FOMP) algorithm of [2] is used, but we stress that other vector selection schemes also can work within the dictionary design framework.

### III. DESIGN EXAMPLES AND DISCUSSION

The framework presented can be used for optimizing virtually any choice of expansion configuration. In this section, we focus on the following established expansions:

- Block transforms.
- frames;
- uniform FIR filterbanks;
- wavelets.

For each case, we will pick a representative expansion and demonstrate that significant improvements are possible using the optimization technique presented earlier. The improvement is measured in terms of increased approximation capability for a given choice of  $S$ , which is the sparseness factor. To emulate

a low bit rate situation, the value used for  $S$  always belong to the set  $\{1/4, 1/6, 1/8, 1/10\}$ , which means that on average only one expansion coefficient out of four, six, eight, or ten signal samples is kept.

In all experiments, the signal  $x(n)$  to be represented is a realization of a Gaussian AR(1) process with correlation factor  $\rho = 0.95$  and with unit variance. The training signal length should be large in order to ensure proper generalization. This means that the optimized signal expansion should give similar approximation results when used on a new realization of the AR(1) process. Here, a signal length of 204 800 samples proved to be sufficient.

#### A. Improving the Karhunen–Loève Transform (KLT)

Although widely considered optimal for compact representation of the chosen Gaussian AR(1) process, it should be emphasized that the KLT only guarantees minimum distortion when used in conjunction with *linear* approximations, i.e., when always retaining the *first* coefficients of a transform block. Optimality is not ensured when the scheme for picking coefficients is based on importance rather than position [20].

Using the  $N = 16$  point KLT as a starting point for the optimization scheme, the resulting waveform dictionary  $\mathcal{F}$  is as shown in the upper left corner of Fig. 1. The goal is now to modify the basis vectors of the KLT in such a manner that the approximating power is increased. In this experiment, the sparseness factor is  $S = 1/4$ , and in each optimization iteration, the coefficients are found using vector selection on a segment of size  $16 \times N = 256$  samples, meaning that 64 out of 256 transform vectors will be selected. This corresponds to a bit allocation situation where a pool of bits are distributed among the coefficients belonging to 16 adjacent transform blocks.

Fig. 4 depicts the obtained training curve showing the approximation error as a function of the number of training iterations. The distortion in iteration 1 is that obtained using the KLT. As explained in Section II-A, each iteration involves two parts: First, find an approximation using vector selection and then compute the best possible waveform dictionary for that set of coefficients. For about 30 iterations, the improvement is monotonous, and then there is a slight increase. The increase in the approximation distortion is a result of the suboptimal vector selection.

The dotted, horizontal line in the plot shows the obtained approximation distortion when retaining the *optimal* set of coefficients/vectors of the KLT. The orthogonal properties of the KLT facilitates easy computation of the best coefficients using the analysis transform. By simply selecting the largest coefficients, the optimal selection is done. In iteration number 1, where the KLT is used as initial transform, we observe that the distortion increase due to suboptimal vector selection is marginal. Furthermore, using the optimized expansion in conjunction with suboptimal vector selection the distortion is reduced by about 15%.

When viewing a traditional transform coder as a subband coder, it is well known that the subband channel responses cover the whole frequency range in a uniform manner [21]. In Fig. 5, we have plotted the amplitude responses for the subband interpretation of a transform coder when using the transform vectors designed above. We observe that the channel responses for

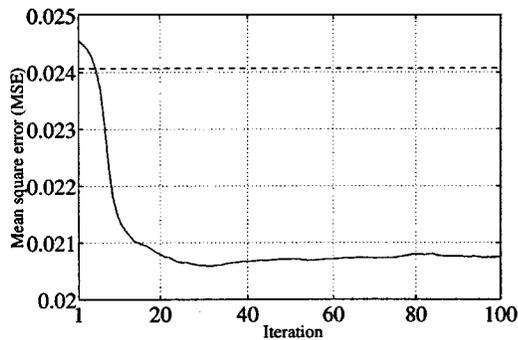


Fig. 4. Training curve using the 16-point KLT as the initial transform and keeping one out of four expansion coefficients. The input signal is a realization of an AR(1),  $\rho = 0.95$  Gaussian process. The mean-square approximation error (MSE) is plotted as a function of the iteration number.

the low-frequency range is very similar to what we would obtain with a DCT or a KLT but that there are no channels in the very-high-frequency range. Given the high degree of sparseness we are aiming at here, this is logical. Rather than computing both high- and low-frequency transform coefficients and then throwing away the high-frequency coefficients, as is often done in classical low bit rate transform coders, we end up with a flexible set of transform vectors, corresponding to overlapping frequency ranges and designed for the best possible approximation of the signal at hand.

We would like to stress that the test of this section is rather harsh. For example, there are no nonstationarities in the examples above to be exploited, but still, the KLT is outperformed. In fact, a similar, or even better performance increase than the one reported above can be realized on real-life nonstationary signals, such as images, speech, and electrocardiogram (ECG) signals. For the case of ECG signals, this was demonstrated in [22].

In a data compression setup, not only the (quantized) coefficients, but also their *position*, would have to be transmitted. Due to the lowpass characteristic of all the basis vectors of the optimized transform, we may anticipate a more uniform distribution of positions than in the KLT case. If entropy coding is used for bit-efficient representation of the position information, the optimized transform will give a higher bit rate for the *position* information. This is not the case if a fixed bit rate coding scheme is used. An investigation into these matters is a topic of future research.

### B. Optimizing Other Configurations

The previous section investigated the optimization of a block transform with the KLT as starting point. The obtained transform was optimized for a sparseness factor of  $S = 1/4$ . Using a different sparseness factor, a different transform would result.

In this section, we focus on approximation capability as a function of the sparseness factor. In the experiments, we use four different signal expansions as starting points:

*Transform:* The  $N = 16$  point KLT, as in Section III-A.

*Frame:* We construct a  $16 \times 32$  frame by merging the basis vectors from an  $N = 16$  point discrete-cosine-transform (DCT) with an  $N = 16$  point Haar transform. This frame was used in [23]–[25] in a data representation setup. The first

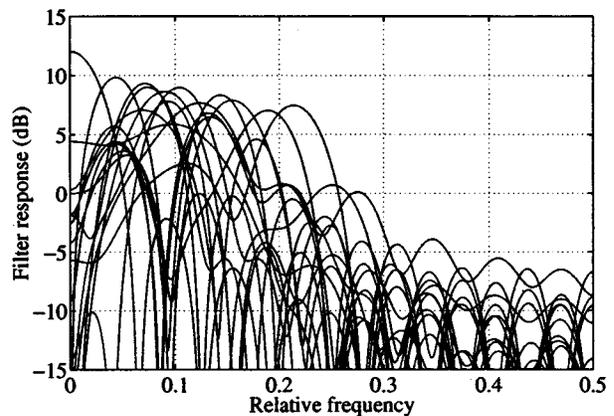


Fig. 5. Frequency responses of the  $N = 16$  channels corresponding to a filterbank interpretation of the optimized transform matrix.

basis vector of these two transforms are identical, and we substitute one of them with a unit-length random vector.

*Wavelet:* We constructed a four-stage dyadic decomposition using the biorthogonal 7/9 wavelet due to Cohen *et al.* [26]. Synthesis of reconstructed signals was done using the nine-tap lowpass filter and the seven-tap highpass filter.

*Uniform FIR Filterbank/LOT:* A 32-tap LOT with  $N = 16$  channels is constructed. This filterbank is optimized for an AR(1)  $\rho = 0.95$  process using an eigenvalue formulation [14].

The exact layout of the associated waveform dictionary  $\mathcal{F}$  for each of these configurations is illustrated in Fig. 1.

Using the same AR(1) training signal as in the previous section, each of the four configurations described above was optimized using the original transform/filterbank coefficients as starting point for the iterations. Each configuration was optimized for  $S = 1/4, 1/6, 1/8,$  and  $1/10$ . As the training curve was not always monotonic, the resulting decomposition was picked as the best one resulting from 100 iterations, except for the frame optimization, where 200 iterations were used due to slower convergence. The approximation capabilities of the obtained decompositions are plotted in Fig. 6 and compared with that of the original decompositions. All plots show the resulting approximation error when using a *new* realization of the AR(1) process. In the cases where the initial decompositions are unitary (KLT and LOT), optimal vector selection is used for the reference curves by retaining the largest coefficients generated using the corresponding analysis decomposition.

For all the chosen configurations, the optimized signal expansions outperform their original counterparts to varying degrees. We observe that the improvement due to optimization increases with the degree of sparseness. The main reason for this is as explained in Section III-A. For very sparse representations, the available expansions vectors should form a flexible set of low-pass vectors. The basis vectors of the KLT corresponding to the higher band channels are hardly used when  $S$  is small. In addition, the use of a greedy vector selection algorithm also contributes to this trend because the vector selection will be closer to optimality when few vectors are selected.

For the three critically sampled configurations, the reduction in approximation error is roughly 30% when one out of ten vec-

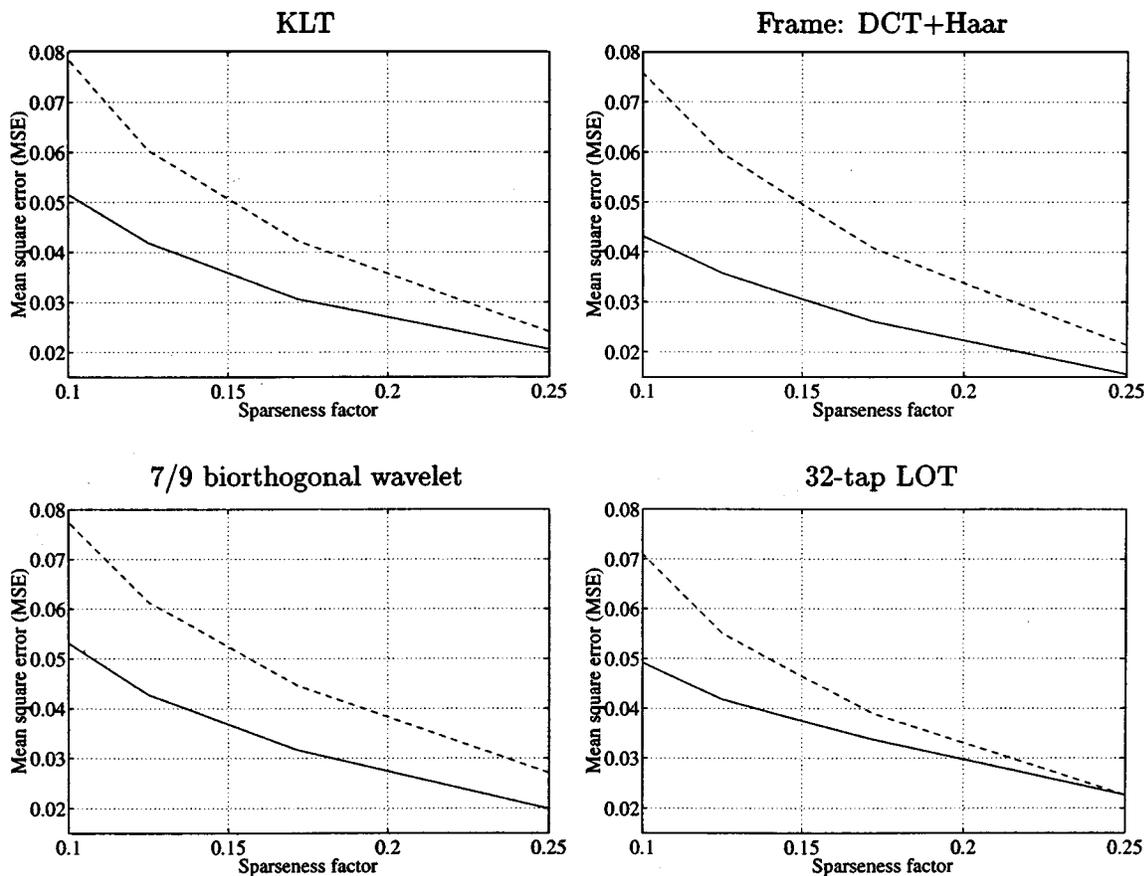


Fig. 6. Comparing the approximation capabilities of established signal expansions (dotted line) with identical configuration where the expansion vectors have been optimized (solid line).

tors are selected, whereas in the frame case, the reduction is about 40%. In addition, we observe that for each sparseness factor, the optimized frame attains the lowest approximation error of the four configurations tested.

#### IV. CONCLUSIONS

In this paper, we have demonstrated that the approximating power of established signal decompositions like the KLT, the original LOT, and the 7/9 biorthogonal wavelet can be improved upon if the objective is to find an expansion of the signal with a given number of terms. This is achieved by giving up the analysis–synthesis paradigm in favor of a simple synthesis approach, where signal expansions are constructed using vector selection algorithms rather than using a corresponding analysis decomposition.

In addition to optimizing signal expansions where the positioning of the expansions vectors are set in order to mimic transforms, filterbanks, or wavelets, the presented algorithm can be used to optimize the approximation capabilities of virtually any kind of signal expansion, including over or undercomplete expansions. We feel that overcomplete expansions show promise for compact representation. This was demonstrated in the experiments where the optimized (overcomplete) frame outperformed the optimized transform in terms of approximation error. However, in a complete data compression setup, the improved approximation capability will have to be paid for by increased side

information giving the positions of the coefficients used. This issue is a topic for further research.

The optimization algorithm is inspired by the GLA used for the design of vector quantizers and involves a training set embedding the properties of the class of signals under consideration. Because the vector selection algorithm is not optimal, the algorithm does not guarantee monotonous decrease of the MSE. In fact, the convergence properties of our design algorithm are not fully understood and is the subject of further investigations.

#### REFERENCES

- [1] B. D. Rao, "Signal processing with the sparseness constraint," in *Proc. ICASSP*, Seattle, WA, May 1998, pp. 1861–1864.
- [2] M. Gharavi-Alkhansari and T. S. Huang, "A fast orthogonal matching pursuit algorithm," in *Int. Conf. Acoust. Speech Signal Process.*, Seattle, WA, May 1998, pp. 1389–1392.
- [3] K. Engan, S. O. Aase, and J. H. Husøy, "Multi-frame compression: Theory and design," *Signal Process.*, vol. 80, pp. 2121–2140, Oct. 2000.
- [4] K. Engan, "Frame based signal representation and compression," Ph.D. dissertation, Norges teknisk-naturvitenskapelige universitet (NTNU)/Høgskolen i Stavanger, Stavanger, Norway, 2000.
- [5] V. K. Goyal, "Quantized overcomplete expansions: Analysis, synthesis and algorithms, tech. rep.," Electron. Res. Lab., memo. UCB/ERL M95/97, 1995.
- [6] V. K. Goyal, M. Vetterli, and N. T. Thao, "Quantization of overcomplete expansions," in *Proc. IEEE Data Compression Conf.*, Snowbird, UT, Mar. 1995, pp. 13–22.
- [7] —, "Quantized overcomplete expansions in RN: Analysis, synthesis, and algorithms," *IEEE Trans. Inform. Theory*, vol. 44, pp. 16–31, Jan. 1998.

- [8] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed in V1," *Vis. Res.*, vol. 37, pp. 3311–3325, 1997.
- [9] M. S. Lewicki and T. J. Sejnowski, "Learning overcomplete representations," *Neural Comput.*, vol. 12, pp. 337–365, Feb. 2000.
- [10] K. Engan, S. O. Aase, and J. H. Husøy, "Designing frames for matching pursuit algorithms," in *Proc. ICASSP*, Seattle, WA, May 1998, pp. 1817–1820.
- [11] ———, "Method of optimal directions for frame design," in *Proc. ICASSP*, Phoenix, AZ, Mar. 1999, pp. 2443–2446.
- [12] A. Gersho, *Vector Quantization and Signal Compression*. Boston, MA: Kluwer, 1992.
- [13] B. A. Ramstad, S. O. Aase, and J. H. Husøy, *Subband Compression of Images—Principles and Examples*. Amsterdam, The Netherlands: Elsevier, 1995.
- [14] H. S. Malvar and D. H. Staelin, "The LOT: Transform coding of images without blocking effects," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 37, pp. 553–559, Apr. 1989.
- [15] H. Malvar, *Signal Processing with Lapped Transforms*. Norwell, MA: Artech House, 1992.
- [16] B. K. Natarajan, "Sparse approximate solutions to linear systems," *SIAM J. Comput.*, vol. 24, pp. 227–234, Apr. 1995.
- [17] S. S. Chen, "Basis Pursuit," Ph.D. dissertation, Stanford University, 1995.
- [18] G. Davis, "Adaptive nonlinear approximations," Ph.D. dissertation, New York Univ., New York, 1994.
- [19] I. F. Gorodnitsky and B. D. Rao, "Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm," *IEEE Trans. Signal Processing*, vol. 45, pp. 600–616, Mar. 1997.
- [20] S. Mallat and F. Falzon, "Understanding image transform codes," in *Proc. SPIE Aerosp. Conf.*, Orlando, FL, Apr. 1997.
- [21] A. N. Akansu and R. A. Haddad, *Multiresolution Signal Decomposition*. San Diego, CA: Academic, 1992.
- [22] J. H. Husøy, S. O. Aase, K. Skretting, and K. Engan, "Design of general block oriented expansions for efficient signal representation," in *Proc. ISCAS*, Orlando, FL, June 1999, pp. III-9–III-12.
- [23] W. Mikhael and A. Ramaswamy, "Application of multitransforms for lossy image representation," *IEEE Trans. Circuits Syst. II*, vol. 41, pp. 431–434, June 1994.
- [24] A. Berg and W. Mikhael, "Signal representation using adaptive parallel mixed transform techniques," in *Proc. 38th IEEE Midwest Symp. Circuits Syst.*, Aug. 1995.
- [25] W. Mikhael and A. Berg, "Image representation using nonorthogonal basis images with adaptive weight optimization," *IEEE Signal Processing Lett.*, vol. 3, pp. 165–167, June 1996.
- [26] A. Cohen, I. Daubechies, and J. C. Feauveau, "Biorthogonal bases for compactly supported wavelets," *Commun. Pure Appl. Math.*, 1992.

**Sven Ole Aase** was born in Stavanger, Norway, in 1965. He received the M.Sc. and Ph.D. degrees in 1989 and 1993, respectively, both from The Norwegian Institute of Technology, Trondheim, Norway.

He is currently a professor with the Department of Electrical and Computer Engineering, Stavanger University College, Stavanger, Norway. His research interests include signal compression and representation, filterbank optimization, adaptive filters, and medical applications of digital signal processing.

**John Håkon Husøy** was born in Toronto, ON, Canada, in 1956. He received the M.Sc. and Ph.D. in electrical engineering in 1981 and 1991, respectively, from the Norwegian Institute of Technology, University of Trondheim, Trondheim, Norway.

He has been involved in hardware and software development in various positions in several companies. He is currently a professor with the Department of Electrical and Computer Engineering, Stavanger University College, Stavanger, Norway. His research interests include image compression, digital filtering, bioelectrical signal processing, adaptive algorithms, and image analysis.

**Karl Skretting** was born in Naerbø, Norway, in 1962. He received the B.Sc. degree in 1985 from Stavanger University College (SUC), Stavanger, Norway. He studied signal processing at Stavanger University College in 1996 and received the M.Sc. degree from the Department of Electrical and Computer Engineering in 1998. Currently, he is a research fellow and is pursuing the Ph.D. degree at SUC with the Signal Processing Group.

His research interests include signal representation and data compression.

**Kjersti Engan** was born in Bergen, Norway, in 1971. She received the B.Sc. degree in electrical engineering from Bergen University College in 1994 and the M.Sc. and Ph.D. degrees, both in electrical engineering, from Stavanger University College (SUC), Stavanger, Norway, in 1996 and 2000, respectively.

She is currently an Associate Professor with the Department of Electrical and Computer Engineering, at SUC. Her research interests include signal and image representation and compression.

Dr. Engan is a member the Norwegian Signal Processing Society.