

Discrete Lagrangian Methods for Optimizing the Design of Multiplierless QMF Banks

Benjamin W. Wah, *Fellow, IEEE*, Yi Shang, *Member, IEEE*, and Zhe Wu

Abstract—In this paper, we present a new discrete Lagrangian method for designing multiplierless quadrature mirror filter banks. The filter coefficients in these filter banks are in powers-of-two, where numbers are represented as sums or differences of powers of two (also called canonical signed digit representation), and multiplications are carried out as additions, subtractions, and shifts. We formulate the design problem as a nonlinear discrete constrained optimization problem, using reconstruction error as the objective, and stopband and passband energies, stopband and passband ripples, and transition bandwidth as constraints. Using the performance of the best existing designs as constraints, we search for designs that improve over the best existing designs with respect to all the performance metrics. We propose a new discrete Lagrangian method for finding good designs and study methods to improve the convergence speed of Lagrangian methods without affecting their solution quality. This is done by adjusting dynamically the relative weights between the objective and the Lagrangian part. We show that our method can find designs that improve over Johnston's benchmark designs using a maximum of three to six ONE bits in each filter coefficient instead of using floating-point representations. Our approach is general and is applicable to the design of other types of multiplierless filter banks.

Index Terms—Adaptive weighted Lagrangian search, canonical signed digit representation, discrete Lagrangian formulation, global search, multiplierless filter banks, quadrature mirror filtering, simulated annealing.

I. INTRODUCTION

DIGITAL filter banks have been applied in many engineering fields. Fig. 1 summarizes the various design objectives for measuring quality. In general, filter-bank design problems are multiobjective, continuous, nonlinear optimization problems.

Algorithms for designing filter banks are either optimization based or nonoptimization based. In optimization-based methods, a design problem is formulated as a multiobjective nonlinear optimization problem [24], whose form may be application- and filter-dependent. The problem is then converted into a single-objective optimization problem and solved by existing optimization methods, such as gradient-descent, Lagrange-multiplier, quasi-Newton, simulated-annealing, and

Manuscript received November 6, 1997; revised May 20, 1999. This research was supported by the National Science Foundation under Grant MIP 96-32316 and by a gift from Rockwell International. This paper was presented in part at the 1997 IEEE International Conference on Application-Specific Array Processors. This paper was recommended by Associate Editor J. Astola.

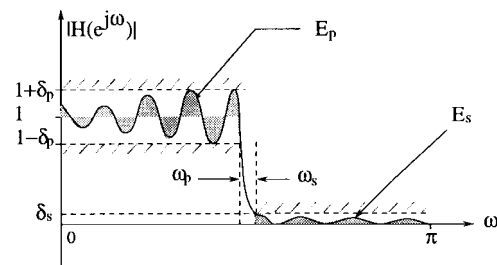
B. W. Wah and Z. Wu are with the Department of Electrical and Computer Engineering and the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA.

Y. Shang is with the Department of Computer Engineering and Science, University of Missouri, Columbia, MI 65211 USA.

Publisher Item Identifier S 1057-7130(99)08038-6.

Filter	Minimization Objectives
Overall	Amplitude distortion
Filter	Aliasing distortion
Bank	Phase distortion
Single	Stopband ripple (δ_s)
Filter	Passband ripple (δ_p)
	Stopband energy (E_s)
	Passband energy (E_p)
	Transition bandwidth (T_t)

(a)



(b)

Fig. 1. Possible design objectives of filter banks and an illustration of the design objectives of a single low-pass filter. $([0, \omega_p]$ is the passband; $[\omega_s, \pi]$ the stopband; $[\omega_p, \omega_s]$, the transition band.)

genetics-based methods [8], [10]. On the other hand, filter bank-design problems have been solved by nonoptimization algorithms, which include spectral factorization [12], [25] and heuristic methods (as in infinite-impulse response (IIR)-filter design). These methods generally do not continue to find better designs once a suboptimal design has been found [25].

In this paper, we study discrete Lagrangian and global-search methods for designing multiplierless quadrature mirror filter (QMF) banks. These filter banks are an important class of filter banks that have been studied extensively. In a two-band QMF bank, the reconstructed signal is

$$\hat{X}(z) = \frac{1}{2}[H_0(z)F_0(z) + H_1(z)F_1(z)]X(z) + \frac{1}{2}[H_0(-z)F_0(z) + H_1(-z)F_1(z)]X(-z) \quad (1)$$

where $X(z)$ is the original signal, and $H_i(z)$ and $F_i(z)$ are, respectively, the response of the analysis and synthesis filters. To perfectly reconstruct the original signal based on \hat{X} , we have to eliminate aliasing, amplitude, and phase distortions. QMF banks with finite-impulse response (FIR) filters implement perfect reconstruction by setting $F_0(z) = H_1(-z)$, $F_1(z) = -H_0(-z)$, and $H_1(z) = H_0(-z)$, leading to a filter bank with one prototype filter $H_0(z)$, linear phase, and no aliasing distortions.

Traditional FIR filters in QMF banks use real or fixed-point numbers as filter coefficients. Multiplications of such

long floating-point numbers generally limit the speed of FIR filtering. To overcome this limitation, *multiplierless (powers-of-two (PO2))* filters have been proposed. These filters use filter coefficients that have only a few bits that are ones. When multiplying a filter input (multiplicand) with one such coefficient (multiplier), the product can be found by adding and shifting the multiplicand a number of times corresponding to the number of ONE bits in the multiplier. For example, the multiplication of y by 0100001001 can be written as the sum of three terms, $y \cdot 2^8 + y \cdot 2^3 + y \cdot 2^0$, each of which can be obtained by shifting y . A limited sequence of shifts and adds are usually much faster than full multiplications. Without using full multiplications, each filter tap takes less area to implement in VLSI, and more filter taps can be accommodated in a given area to implement filter banks of higher performance.

The frequency response of a PO2 filter $H(z)$ is

$$H(z) = \sum_{i=0}^{\gamma-1} x_i z^{-i} = \sum_{i=0}^{\gamma-1} \left(\sum_{j=0}^{d-1} e_{i,j} 2^j \right) z^{-i}$$

where

$$\sum_{j=0}^{d-1} |e_{i,j}| \leq \text{for all } i, \quad e_{i,j} = -1, 0, 1. \quad (2)$$

Here, γ is the length of the PO2 filter, l is the maximum number of ONE bits used in each coefficient, and d is the number of bits in each coefficient.

The design of multiplierless filters has been solved by integer programming that optimizes filter coefficients with restricted values of PO2. Other techniques used include combinatorial search [17], simulated annealing [2], genetic algorithms [18], linear programming [11], and continuous Lagrange-multiplier methods in combination with a tree search [20].

In this paper, we present a discrete Lagrange-multiplier search for designing multiplierless QMF banks. In Section II, we formulate the design problem as a single-objective constrained optimization problem. Section III summarizes the principles behind discrete Lagrangian methods. In Section IV, we present our discrete Lagrangian method, 1998 version (DLM-98), that finds saddle points in discrete space and examines the issues related to the implementation of DLM-98 to design multiplierless filter banks. Finally, Section V presents experimental results, and conclusions are drawn in Section VI.

II. PROBLEM FORMULATION

The design of QMF banks can be formulated as a multi-objective unconstrained optimization problem or as a single-objective constrained optimization problem.

A. Multi-Objective Unconstrained Formulation

In a multiobjective formulation, the goals can be to:

- 1) minimize the amplitude distortion (reconstruction error) of the overall filter bank;
- 2) optimize the individual performance measures of the prototype filter $H_0(z)$.

One possible formulation using a subset of the measures in Fig. 1 is as follows:¹

Minimize E_r and E_s where

$$E_r = \int_{\omega=0}^{\pi/2} \cdot \left(|H_0(e^{j\omega})|^2 + |H_0(e^{j(\omega-\pi)})|^2 - 1 \right)^2 d\omega$$

and

$$E_s = \int_{\omega=\omega_s}^{\pi} |H_0(e^{j\omega})|^2 d\omega. \quad (3)$$

Unfortunately, optimal solutions to (3) are not necessarily optimal solutions to the original problem that considers all the performance measures. Often, performance measures not included in the formulation are compromised.

In general, optimal solutions of a multiobjective problem form a *Pareto optimal frontier* such that one solution on this frontier is not dominated by another. One approach to find a point on the frontier is to optimize a weighted sum of all the objectives [3], [6], [10], [16], [24]. This approach has difficulty when frontier points of certain characteristics are desired, such as those with certain transition bandwidth. Different combinations of weights must be tested by trial and error until a desired solution is found. When the desired characteristics are difficult to satisfy, trial and error is not effective in finding feasible designs. Instead, constrained formulations should be used.

B. Single-Objective Constrained Formulation

Another approach to solve a multiobjective problem is to turn all but one objectives into constraints, and define the constraints with respect to a reference design. The specific measures constrained may be application and filter dependent [24].

Constraint-based methods have been applied to design QMF banks in both the frequency [3], [5], [10], [21], [23] and time domains [15], [22]. In the frequency domain, the most often considered objectives are reconstruction error (E_r) and (stopband ripple) (δ_s). As stopband ripples cannot be formulated in closed form, stopband attenuation is used instead (represented as E_s in Fig. 1). In the time domain, Nayebi [15] gave a time-domain formulation with constraints in the frequency domain, and designed filter banks using an iterative time-domain design algorithm.

In this paper, we formulate the design of QMF banks in the most general form as a nonlinear constrained optimization problem using the reconstruction error as the objective, and other measures (stopband ripple, stopband energy, passband ripple, passband energy and transition bandwidth) as constraints

Minimize E_r , subject to

$$E_p \leq \theta_{E_p}, \quad E_s \leq \theta_{E_s}, \quad T_t \leq \theta_{T_t}, \quad \delta_p \leq \theta_{\delta_p}, \quad \delta_s \leq \theta_{\delta_s} \quad (4)$$

where θ_{E_p} , θ_{E_s} , θ_{δ_p} , θ_{δ_s} , and θ_{T_t} are constraint bounds found in the best-known design (with possibly some bounds relaxed or tightened in order to obtain designs of different tradeoffs).

¹Note that in QMF banks, E_r is nonzero. A multirate filter bank that enforces perfect reconstruction ($E_r = 0$) can be formulated as a constrained optimization problem with a goal of minimizing E_s [8], [9].

The goal here is to find filter banks of a finite word length whose performance measures are better than, or equal to, those of the reference design. Since the objective and the constraints are nonlinear, the problem is multimodal with many local minima.

The original optimization problem with inequality constraints (4) can be transformed into an optimization problem with equality constraints as follows:

$$\begin{aligned}
 \text{Minimize } f(x) &= V_{E_r} = \frac{E_r - \theta_{E_r}}{\theta_{E_r}} & (5) \\
 \text{subject to } V_{E_p} &= \max\left(\frac{E_p - \theta_{E_p}}{\theta_{E_p}}, 0\right) = 0 \\
 V_{E_s} &= \max\left(\frac{E_s - \theta_{E_s}}{\theta_{E_s}}, 0\right) = 0 \\
 V_{\delta_p} &= \max\left(\frac{\delta_p - \theta_{\delta_p}}{\theta_{\delta_p}}, 0\right) = 0 \\
 V_{\delta_s} &= \max\left(\frac{\delta_s - \theta_{\delta_s}}{\theta_{\delta_s}}, 0\right) = 0 \\
 V_{T_t} &= \max\left(\frac{T_t - \theta_{T_t}}{\theta_{T_t}}, 0\right) = 0 & (6)
 \end{aligned}$$

where x is a vector of discrete coefficients, θ_{E_r} is the reconstruction error of the best-known design, and all functions have been normalized with respect to the values of the best-known design.

III. LAGRANGIAN FORMULATIONS AND METHODS

In this section, we first summarize past work on Lagrangian formulations and methods for solving continuous constrained optimization problems. We then extend them to discrete constrained optimization problems [19], [26], [29].

A. Continuous Lagrangian Formulations and Methods

Lagrangian methods are classical methods for solving continuous constrained optimization problems [14]. We first review briefly the theory of Lagrange multipliers.

Define a continuous constrained optimization problem as follows:

$$\begin{aligned}
 \min_{x \in E^m} f(x) \\
 \text{subject to } g(x) &\leq 0 \quad x = (x_1, x_2, \dots, x_n) \\
 h(x) &= 0 & (7)
 \end{aligned}$$

where x is a vector of real numbers, $f(x)$ is an objective function, $g(x) = [g_1(x), \dots, g_k(x)]^T$ is a set of k inequality constraints, and $h(x) = [h_1(x), \dots, h_m(x)]^T$ is a set of m equality constraints. Further, $f(x)$, $g(x)$, and $h(x)$, as well as their derivatives, are continuous functions.

Since Lagrangian methods cannot deal directly with inequality constraints, we transform inequality constraint $g_i(x) \leq 0$ into an equality constraint by adding a slack variable $z_i(x)$, transforming it into $p_i(x) = g_i(x) + z_i^2(x) = 0$. The transformation is done in such a way that guarantees the existence of gradients in the x space when $g_i(x) = 0$. The

corresponding *Lagrangian function* is defined as follows:

$$L_c(x, \lambda, \mu) = f(x) + \lambda^T h(x) + \mu^T p(x) + \frac{1}{2} \|p(x)\|^2 \quad (8)$$

where $\lambda = [\lambda_1, \dots, \lambda_m]^T$ and $\mu = [\mu_1, \dots, \mu_k]^T$ are two sets of Lagrange multipliers, and $p(x) = [p_1(x), \dots, p_k(x)]^T$.

To eliminate z_i from (8), we minimize L_c with respect to z_i for a given x . After substituting the result into (8) [14], we have

$$\begin{aligned}
 L_c(x, \lambda, \mu) &= f(x) + \lambda^T h(x) \\
 &+ \frac{1}{2} \sum_{i=1}^k \left[\max(0, \mu_i + g_i(x)) - \mu_i^2 \right]. & (9)
 \end{aligned}$$

Note that the derivation applies to both the continuous and the discrete cases because the differentiation of L_c with respect to z_i is for a fixed x , and z_i is assumed continuous.

According to classical optimization theory [14], all the extrema of (9) that satisfy the constraints and that are regular points are roots of the following set of first-order necessary conditions:

$$\begin{aligned}
 \nabla_x L_c(x, \lambda, \mu) &= 0 \\
 \nabla_\lambda L_c(x, \lambda, \mu) &= 0 \\
 \nabla_\mu L_c(x, \lambda, \mu) &= 0. & (10)
 \end{aligned}$$

These conditions are necessary to guarantee the (local) optimality of the solution to (7).²

There are many ways to find solutions that satisfy (10), including sequential quadratic programming and first-order search methods. The *first-order search method* expresses the search in a dynamic system of ordinary differential equations:

$$\begin{aligned}
 \frac{d}{dt} x(t) &= -\nabla_x L_c(x(t), \lambda(t), \mu(t)) \\
 \frac{d}{dt} \lambda(t) &= \nabla_\lambda L_c(x(t), \lambda(t), \mu(t)) \\
 \frac{d}{dt} \mu(t) &= \nabla_\mu L_c(x(t), \lambda(t), \mu(t)). & (11)
 \end{aligned}$$

They perform *local search* involving simultaneous descents in the original-variable space of x and ascents in the Lagrange-multiplier space of λ and μ . They evolve over time t , and reach a feasible local extremum when they stop at an *equilibrium point* where all gradients are zeros.

B. Discrete Lagrangian Formulations and Methods

For discrete optimization problems, all the variables x_i ($i = 1, 2, \dots, n$) take discrete values (e.g., integers). Little work has been done in applying Lagrangian methods to solve discrete constrained combinatorial optimization problems [7]. The difficulty in traditional Lagrangian methods lies in the lack of a differentiable continuous space to find an equilibrium point. In this subsection, we describe the theory of Lagrange-multiplier methods in discrete space [19], [29].

²There are second-order conditions to guarantee that the extremum found is a local minimum [14].

For nonlinear discrete problems with inequality constraints (similar to that in (7) where x is now a vector of discrete variables), we first transform inequality constraint $g_i(x) \leq 0$ into an equality constraint $\max(g_i(x), 0) = 0$. This transformation does not use a slack variable as in the continuous case because searches in discrete Lagrangian space do not require the existence of gradients in the x space when $g_i(x) = 0$.

After transforming inequality constraints into equality constraints, the resulting discrete Lagrangian function is written as follows:

$$L_d(x, \lambda, \mu) = f(x) + \lambda^T h(x) + \sum_{i=1}^k \mu_i \max(0, g_i(x)) \quad (12)$$

where x is a discrete variable and λ and μ can be continuous.

The discrete Lagrangian function defined in (12) cannot be used to derive similar conditions in (10) because there are no gradients and differentiation in discrete space. Without these concepts, none of the mechanisms of calculus in continuous space is applicable in discrete space.

An understanding of gradients in continuous space shows that they define directions in a small neighborhood in which function values decreases. To this end, we define in discrete space a *direction of maximum potential drop* for Lagrangian function $L_d(x, \lambda, \mu)$ at point x for fixed λ and μ as a vector³ that points from x to a neighborhood point of $x \in \mathcal{N}(x)$ with the minimum L_d value

$$\begin{aligned} \Delta_x L_d(x, \lambda, \mu) &= \vec{v}_x \\ &= y \ominus x \\ &= (y_1 - x_1, y_2 - x_2, \dots, y_n - x_n) \end{aligned}$$

where

$$y \in \mathcal{N}(x) \cup \{x\} \quad (13)$$

$$L_d(y, \lambda, \mu) = \min_{\substack{x' \in \mathcal{N}(x) \\ \cup \{x\}}} L_d(x', \lambda, \mu).$$

Here, \ominus is the vector-subtraction operator for changing x in discrete space to one of its “user-defined” neighborhood points $\mathcal{N}(x)$. Intuitively, \vec{v}_x is a vector pointing from x to y , the point with the minimum L_d value among all neighboring points of x , including x itself. That is, when x itself has the minimum L_d , then $\vec{v}_x = \vec{0}$.

Having defined $\Delta_x L_d(x, \lambda, \mu)$, the direction of maximum potential drop in the x space, we define the concept of *saddle points* in discrete space similar to those in continuous space [14]. A point (x^*, λ^*, μ^*) is a saddle point when

$$L(x^*, \lambda, \mu) \leq L(x^*, \lambda^*, \mu^*) \leq L(x, \lambda^*, \mu^*) \quad (14)$$

for all $x \in \mathcal{N}(x^*)$, all possible λ , and all possible $\mu \geq \mu^*$. Starting from (14), we can prove similar first-order necessary conditions in discrete space that are satisfied by all saddle points [29]

$$\begin{aligned} \Delta_x L_d(x, \lambda, \mu) &= 0 \\ \nabla_\lambda L_d(x, \lambda, \mu) &= 0 \\ \nabla_\mu L_d(x, \lambda, \mu) &= 0. \end{aligned} \quad (15)$$

³We assume that points in the x space are represented as vectors without explicitly denoting them using the vector notation, whereas λ and μ are scalars.

Note that the notation in the first condition defines the direction of maximum potential drop of L_d in discrete space of x for fixed λ and μ , whereas the differentiations in the last two conditions are in continuous space of λ and μ for fixed x . For brevity, the proofs showing the correctness of these conditions are omitted here [29].

The first-order necessary conditions in (15) lead to the following first-order search method in discrete space. Here, we seek discrete equilibrium points similar to those of continuous problems. The following equations are discrete approximations to implement the first-order conditions in (15).

General Discrete First-Order Search Method:

$$x(k+1) = x(k) \oplus \Delta_x L_d(x(k), \lambda(k), \mu(k)) \quad (16)$$

$$\lambda(k+1) = \lambda(k) + c_1 h(x(k)) \quad (17)$$

$$\mu(k+1) = \mu(k) + c_2 \max(0, g(x(k))) \quad (18)$$

where \oplus is the vector-addition operator ($x \oplus y = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n)$), and c_1 and c_2 are positive real numbers controlling how fast the Lagrange multipliers change.

It is easy to see that the necessary condition for the discrete first-order method to converge is when $h(x) = 0$ and $g(x) \leq 0$, implying that x is a feasible solution to the original problem. If any of the constraints is not satisfied, then λ and μ on the unsatisfied constraints will continue to evolve. Note that, as in continuous Lagrangian methods, the first-order conditions are only satisfied at saddle points, but does not imply that the time to find a saddle point is finite, even if one exists.

IV. DLM-98: AN IMPLEMENTATION OF DISCRETE FIRST-ORDER METHOD

Based on (5) and (6), the discrete Lagrangian function for optimizing PO2 filter banks is

$$L_d(x, \lambda, \mu) = f(x) + \sum_{\substack{i \in \{E_p, E_s, \\ \delta_p, \delta_s, T_i\}}} \mu_i \cdot V_i \quad (19)$$

where x is a vector of coefficients, each of which is in canonical signed digit (CSD) form of the sum of several signed binary bits, such as $2^{-1} + 2^{-3} - 2^{-6}$. Since we have only equality constraints transformed from inequality constraints, we use μ as our Lagrange multipliers in the following discussion.

Fig. 2 shows an implementation of the discrete first-order method (16) and (18) for designing PO2 filter banks formulated as nonlinear discrete constrained minimization problems. The procedure shows several aspects that can be tuned in order to improve its performance.

- *Starting Points (Line 3):* we choose a starting point based on a discrete approximation of an existing QMF bank with real coefficients (Section IV-A).
- *Initial Lagrange-Multiplier Values (Line 4):* we initialize all Lagrange multipliers to zero in order to allow our results to be reproduced easily. An optimal initial setting is difficult because it depends on the amount of constraint violation.
- *Time Constraint (Lines 2 and 6):* this limits the number of iterations through the loop.

procedure *DLM-98*

1. set c (positive real constant for controlling the speed of change of Lagrange multipliers);
2. set i_{max} (maximum number of iterations);
3. set starting point x ;
4. set initial value of μ (set to 0 in the experiments);
5. **if** using dynamic weight adaptation **then** *weight_initialization*;
6. **while** search has not converged **and** number of iterations $< i_{max}$ **do** {
7. update x to x' only if this will result in $L_d(x', \mu) < L_d(x, \mu)$;
8. **if** condition for updating μ is satisfied **then** $\mu_i \leftarrow \mu_i + c \cdot \max(0, g_i)$;
9. **if** using dynamic weight adaptation **then** *dynamic_weight_adaptation* }

Fig. 2. An implementation of discrete first-order method for designing PO2 filter banks. (The initial values of parameters are indicated here unless specified otherwise in the text.)

- *Updating x (Line 7):* here, we evaluate all possible neighboring points of x in order to find improvements in its Lagrangian value (Section IV-B).
- *Updating μ (Lines 1 and 8):* the Lagrange multipliers are updated when the search reaches a local minimum in the objective space. We do not update the multipliers more frequently due to instability of the trajectory. The amount of update is controlled by an application-dependent constant c and other filter-related parameters (Section IV-C).
- *Dynamic Weight Adaptation (Lines 5 and 9):* weight adaptation adjusts the weight between the objective and the constraints in order to adjust their relative importance and to improve convergence (Section IV-D).

A. Generating a Starting Point

There are two alternatives to select a starting point (Line 3 in Fig. 2): using the parameters of an existing PO2 QMF bank, or using a discrete approximation of an existing QMF bank with real coefficients. The first alternative is not always possible because not many such filter banks are available in the literature. In this section, we discuss the second alternative.

In the second approach, we first transform the real coefficients of the best-known design to PO2 forms using a CSD representation. Given a real coefficient and b , the maximum number of ONE bits to represent the coefficient, we apply Booth’s algorithm [1] to represent consecutive 1’s using two ONE bits and then truncate the least-significant bits of the coefficients. This approach generally allows a number to be represented in a few ONE bits. As an example, consider a binary fixed-point number 0.10011101100. After applying Booth’s algorithm and truncation, we can represent the number in two ONE bits

$$\begin{array}{lcl}
 0.10011101100 & \xrightarrow{\text{Booth's Algorithm}} & \\
 0.101000\bar{1}0\bar{1}00 & \xrightarrow{\text{Truncation}} & 2^{-1} + 2^{-3}.
 \end{array}$$

Previous work [4], [13], [17], shows that scaling has a significant impact on the optimization of coefficients in PO2 filters. That is, if each coefficient is scaled properly before the search starts (based on a heuristic objective), the quality of the final design can be improved significantly. In our case, the performance of a PO2 filter obtained by truncating its real

TABLE I
COMPARISON OF A PO2 FILTER BANK OBTAINED BY TRUNCATING THE REAL COEFFICIENTS OF JOHNSTON’S 32e QMF BANK [10] TO THREE BITS AND A SIMILAR PO2 FILTER BANK, WHOSE COEFFICIENTS WERE SCALED BY 0.5565 BEFORE TRUNCATION

Performance Metrics	E_r	E_p	E_s	δ_p	δ_s	T_t
Filter bank A with Truncated Coefficients	6.93	9.61	1.09	1.89	1.05	1.00
Filter bank B with Scaling and Truncation	0.99	1.08	0.96	1.20	0.98	0.99

(Performance has been normalized with respect to the performance of the original filter bank).

coefficients to a fixed maximum number of ONE bits is not as good as one whose real coefficients were first multiplied by a scaling factor. We illustrate this observation in the following example.

Consider Johnston’s 32e filter bank [10] as a starting point. Table I shows the metrics of two PO2 filters: Filter Bank A was obtained by truncating each of the original coefficients to a maximum of three ONE bits, whereas Filter Bank B was obtained by multiplying each of the coefficients by 0.5565 before truncation. Filter Bank B performs better and is almost as good as the original design with real coefficients. In fact, a design that is better than Johnston’s 32e design can be obtained by using Filter B as a starting point, but no better designs were found using Filter A. This example illustrates that multiplying the filter coefficients by a scaling factor changes the bit patterns of the coefficients, which can improve the quality of the starting point when the coefficients are truncated.

Experiments also show that it is possible to find good designs without requiring the PO2 coefficients to have the same degree of precision as that of continuous coefficients. For instance, in our experiments, we restrict the minimum exponent of the ONE bits in each coefficient (in the range $[-1, 1]$) to be -22 , even though the real coefficients have a minimum exponent of -31 .

To find the best scaling factor, we enumerate over different scaling constants and scale all the coefficients by a common constant before the search begins. Fig. 3 shows a simple but effective algorithm to find the proper scaling factor to be multiplied before the coefficients are truncated. We evaluate the quality of the resulting starting point by a weighted sum of its performance metrics. Since, under most circumstances, the constraint on transition bandwidth is more difficult to satisfy, we give it a weight of 100 and a weight of 1 for the other four metrics. Note that our objective in finding a good scaling

procedure *find_scaling_factor*

1. $LeastSum = +\infty$;
2. **for** $ScaleFactor := 0.5000$ **to** 1.0 **step** 0.0001 **do** {
3. Multiply each filter coefficient by $ScaleFactor$;
4. Get the PO2 form of the scaled coefficients;
5. Compute the weighted sum of constraint violation: $sum := \sum_{i=1}^5 w_i \cdot g_i$;
6. **if** ($sum < LeastSum$) **then** { $LeastSum := sum$; $BestScale := ScaleFactor$ }
7. **return** $BestScale$

Fig. 3. Algorithm for finding the best scaling factor, where w_i is the weight of constraint i .

factor is different from that in the previous work [4], [13], [17]. Further, note that the filter output in the final design will need to be divided by the same scaling factor.

Experimental results show that the algorithm in Fig. 3 works fast and can complete in a few minutes, and that the scaling factors chosen are reasonable and suitable. It is important to point out that scaling does not help when the number of ONE bits allowed to represent each coefficient is large. For instance, when the maximum number of ONE bits allowed is larger than six, the performance of all the filters is nearly the same for all scaling factors.

As an illustration, consider the design of a PO2 QMF bank [28] based on Johnston's 32d design [10] as our constraints. Assuming a minimum exponent of -22 in each ONE bit, we enumerate and find the best scaling factor for all the coefficients to be 0.9474.

B. Updating x

The value of x is updated in Line 7 in Fig. 2. There are two ways in which x can be updated: greedy update and hill climbing. In greedy updates, the update of x leading to the maximum improvement of $L_d(x, \mu)$ is found before an update is made. This approach is very time consuming and may not lead to the best filter bank when DLM-98 stops. On the other hand, in hill climbing, x is updated as soon as an improvement in $L_d(x, \mu)$ is found. This approach is efficient and generally leads to good designs. For this reason, we use hill climbing as our update strategy.

We process all the bits of all the coefficients in a round-robin manner. Suppose γ is the filter length, l is maximum number of ONE bits that can be used for each coefficient, and the i th coefficient is composed of l elements $b_{i,1}, b_{i,2}, \dots, b_{i,l}$. We process the elements in the following order repetitively:

$$b_{1,1}, b_{1,2}, \dots, b_{1,l}, b_{2,1}, \dots, b_{\gamma,1}, \dots, b_{\gamma,l}.$$

For each element $b_{i,j}$, we perturb it to be $b'_{i,j}$ that differs from $b_{i,j}$ by either the sign or the exponent or both, while maintaining $b'_{i,j}$ to be not the same in exponent as another element of the i th coefficient. Using $b_{i,1}, \dots, b_{i,j-1}, b'_{i,j}, \dots, b_{i,l}$ while keeping other coefficients the same, we compute the new value $L_d(x', \mu)$ and accept the change if $L_d(x', \mu) < L_d(x, \mu)$.

C. Updating μ

Lines 1 and 8 in Fig. 2 are related to the condition when μ should be updated. In traditional Lagrangian methods on

continuous variables, μ is updated in every iteration. This approach does not work in DLM-98 because if μ were updated after each update of x , then the search behaves like random probing and restarts from a new starting point even before a local minimum is reached. For this reason, μ for violated constraints should be updated less frequently, only when no further improvement in $L_d(x, \mu)$ can be made in Line 7 of DLM-98 for all the bits in all the coefficients. This is the approach we have taken in solving satisfiability problems [19], [26], [29]. However, we have found that more frequent updates of μ may lead to better PO2 filters. In our implementation, we update μ every time three coefficients have been processed. Since μ is updated before all the filter coefficients have been perturbed, the guidance provided by μ may not be exact.

When updating μ before the search reaches a local minimum of $L_d(x, \mu)$, we set c in Line 8 of Fig. 2 to be a normalized value as follows:

$$c = \frac{\theta_{\text{speed}}}{\max_{i=1}^n g_i} \quad (20)$$

where θ_{speed} is a real constant for controlling the speed of increasing μ . Experimentally, we have determined θ_{speed} to be 0.6818.

When the search reaches a local minimum of $L_d(x, \mu)$, perturbing any single bit in any coefficient will result in no improvement of $L_d(x, \mu)$. At this point, we need to update μ differently in order to bring the search out of the local minimum. This is done by choosing a proper value of c in Line 8 of DLM-98. If μ is increased too fast, then the search will restart from a random starting point. On the other hand, if μ is increased too slowly, then the trajectory will remain in the current local minimum, and updates of x in the next iteration of DLM-98 will bring the search to the same local minimum! Hence, we like to set c so that it will bring the search out of the current local minimum in one step, and local descents in the next iteration will head to an adjacent local minimum. This means that, after μ has been changed to μ' , there exists x' in $\mathcal{N}(x)$, such that

$$L_d(x, \mu) \leq L_d(x', \mu) \quad \text{and} \quad L_d(x', \mu') < L_d(x, \mu'). \quad (21)$$

Replacing $L_d(x, \mu)$ by $f(x) + \sum_{i=1}^n \mu \cdot \max(0, g_i(x))$ in (21), we get the condition before μ changes

$$f(x) + \sum_{i=1}^n \mu \cdot \max(0, g_i(x)) \leq f(x') + \sum_{i=1}^n \mu \cdot \max(0, g_i(x')) \quad (22)$$

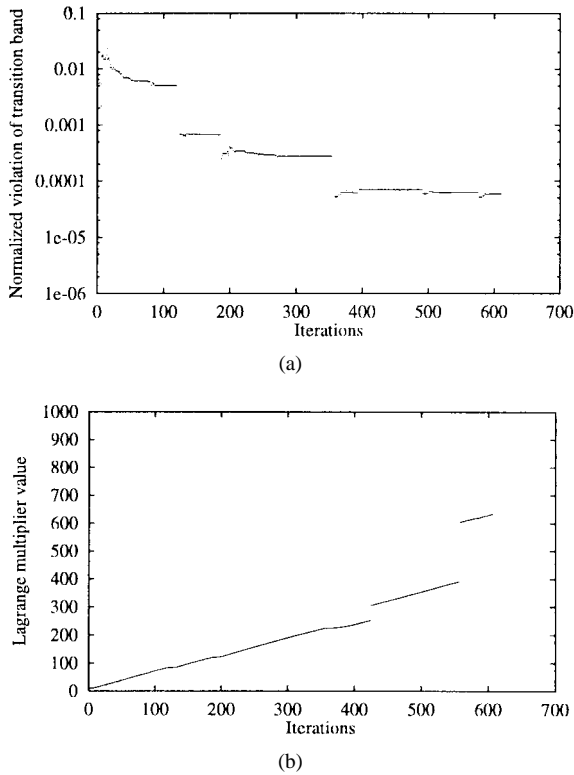


Fig. 4. Performance progress measured during the search of Problem 32e. (a) Violation values of T_t . (b) Corresponding μ_{T_t} .

and that after μ is updated to $\mu'_i = \mu_i + c \cdot \max(0, g_i(x))$

$$f(x) + \sum_{i=1}^n \mu'_i \cdot \max(0, g_i(x)) > f(x') + \sum_{i=1}^n \mu'_i \cdot \max(0, g_i(x')) \quad (23)$$

where $\max(0, g_i(x'))$ is the new violation of the i th constraint at x' . After transformations, we get

$$c > \frac{L_d(x', \mu) - L_d(x, \mu)}{\sum_{i=1}^n \max(0, g_i(x)) \cdot (\max(0, g_i(x)) - \max(0, g_i(x')))} \quad (24)$$

When c is large enough to satisfy (24) for all x' , and μ is increased according to Line 8 of DLM-98, we are assured that there is new x' that will cause L_d to decrease in the next iteration.

As an example, consider in Fig. 4(a) the violation of transition bandwidth T_t in a typical search based on the constraints derived from Johnston's 32e filter bank [10]. Fig. 4(a) shows that the value of the violation on T_t can be extremely small, on the order of 10^{-5} in the later part of the search. For such small violation values, the update of μ_{T_t} using c defined in (20) will result in a large number of iterations before the violation can be overcome. Using c defined in (24) to increase μ_{T_t} , we see in Fig. 4(b) that μ_{T_t} jumps three times when the condition for updating μ was satisfied. These saved at least half of the total search time in order to find the solution.

TABLE II
MULTIPLIERLESS 32d QMF BANKS FOUND BY DLM-98 WITH STATIC AND ADAPTIVE WEIGHTS (THE OBJECTIVE IS THE RECONSTRUCTION ERROR E_r)

Weight w	Static Weight		Adaptive Weight	
	Objective	Time (minutes)	Objective	Time (minutes)
10000.0	-	-	0.998	195.9
1000.0	-	-	0.998	168.3
100.0	-	-	0.885	277.2
10.00	-	-	0.883	119.3
1.00	-	-	0.836	190.8
0.1	0.87	197.1	0.837	161.4
0.01	0.87	115.2	0.87	124.3
0.001	0.87	4.8	0.87	34.5
0.0001	0.88	12.0	0.878	10.8
0.00001	0.9	23.7	0.924	12.0

D. Weighted Discrete First-Order Methods

As discussed in Section III-A, Lagrangian methods rely on ascents in the Lagrange-multiplier space, and descents in the objective space, in order to reach equilibrium. The convergence speed and solution quality, however, depend on the balance between objective $f(x)$ and constraints $h(x)$ and $g(x)$. Although changes in μ lead to different balance between ascents and descents, convergence can be improved by introducing a weight on the objective function. These considerations lead to a new Lagrangian function as follows:

$$L_d^w(x, \mu) = wf(x) + \sum_{i=1}^k [\mu_i \max(0, g_i(x))] \quad (25)$$

where $w > 0$ is a user-controlled weight on the objective. By applying DLM-98 in Fig. 2 on (25) using different w , we observe four possible behaviors of the search trajectory.

- 1) The trajectory converges without oscillations.
- 2) The trajectory gradually reduces in oscillations and eventually converges.
- 3) The trajectory oscillates within some range but never converges.
- 4) The magnitude of oscillations increases, and the trajectory eventually diverges.

Obviously, the first two cases are desirable, and the latter two are not. Moreover, we would like to reduce the amount of oscillations and improve convergence time.

The second and third columns of Table II show the objective-function values of the designs found and the corresponding convergence times of DLM-98 with static weights. DLM-98 does not converge when the static weight w is large. These results demonstrate that the choice of w is critical in controlling both the convergence time and solution quality. There is, however, no effective method for choosing a fixed w except by trial and error.

In the rest of this subsection, we present a strategy to adapt w based on run-time search progress in order to obtain high-quality solutions and short convergence time. This approach is more general than our previous approach [19] that scales the Lagrange multipliers periodically in order to prevent them from growing to be very large when all constraint functions are positive. The Lagrange multiplier of a nonnegative constraint may grow without limit because its value is always nondecreasing according to (18), and a

procedure *weight_initialization*

1. set $w(0)$ (initial weight, set to 0.00001 in the experiments);
2. set N_u (major window for changing w , set to 30 in the experiments);
3. set δ_t (minor window for changing w , set to 5 in the experiments);
4. $j \leftarrow 0$ (number of iterations since last divergence)

procedure *dynamic_weight_adaptation*

5. record useful information for calculating performance;
6. $j \leftarrow j + 1$;
7. **if** ($j \bmod \delta_t = 0$) **then**
8. **if** trajectory diverges **then** { reduce w ; $j \leftarrow 0$ }
9. **if** ($j \bmod N_u = 0$) **then** {
10. compute performance metrics based on data collected;
11. change w when certain conditions are satisfied (see text) }

Fig. 5. Procedures for weight initialization and adaptation in Fig. 2. (The initial values of parameters are indicated here unless specified otherwise in the text.)

Lagrangian space with large Lagrange multipliers is more rugged and more difficult to search. In our previous approach [19], the period between scaling and the scaling factor are application dependent and chosen in an ad hoc fashion. Our current approach adjusts the weight between the objective and the constraints, which is equivalent to scaling the Lagrange multipliers. It is more general because it adjusts the weight according to the convergence behavior of the search.

In general, changing w may speed up or delay convergence before a trajectory reaches an equilibrium point, and may bring the trajectory out of equilibrium after it reaches there. In this section, we design weight-adaptation algorithms to speed up convergence. Strategies to bring a trajectory out of equilibrium by modifying w will be studied in the future.

Fig. 5 outlines the procedures for weight initialization and adaptation. Its basic idea is to first estimate the initial weight $w(0)$ (Line 1), measure the performance of the search trajectory $(x(t), \mu(t))$ periodically, and adapt $w(t)$ to improve convergence time or solution quality.

Let (x_i, μ_i) be the point of the i th iteration, and $v_{\max}(i)$ be its maximum violation over all the k constraints in (25)

$$v_{\max}(i) = \max_{1 \leq j \leq k} \{g_j(x_i), 0\}. \quad (26)$$

To monitor the progress of the search, we divide time into nonoverlapping major windows of size N_u iterations (Line 2), each of which is then divided into minor windows of δ_t iterations (Line 3). We further record some statistics, such as $v_{\max}(i)$ and $f_i(x)$, that will be used to calculate the performance in each minor/major window (Line 5).

At the beginning of a minor window (Line 7), we test whether the trajectory diverges or not (Line 8). Divergence happens when $v_{\max}(i)$ is larger than an extremely large value (say 10^{20}). If it happens, we reduce w , say $w \leftarrow w/10$, and restart the window markers by resetting j to zero.

At the beginning of a major window (Line 9), we compute some metrics to measure the progress of the search relative to that of previous major windows (Line 10). In general, application-specific metrics, such as the number of oscillations of the trajectory, can be used. In our current implementation,

we compute the averages (or medians) of $v_{\max}(i)$ and objective $f_i(x)$ in the u th major window ($u = 1, 2, \dots$) as follows:

$$\bar{v}_u = \frac{1}{N_u} \sum_{i=(u-1)N_u+1}^{uN_u} v_{\max}(i) \quad \text{or} \quad \bar{v}_u = \text{median}_{\substack{(u-1)N_u+1 \\ \leq i \leq uN_u}} \{v_{\max}(i)\} \quad (27)$$

$$\bar{f}_u = \frac{1}{N_u} \sum_{i=(u-1)N_u+1}^{uN_u} f_i(x) \quad \text{or} \quad \bar{f}_u = \text{median}_{\substack{(u-1)N_u+1 \\ \leq i \leq uN_u}} \{f_i(x)\} \quad (28)$$

Based on these measurements, we adjust w accordingly (Line 11). Note that when comparing values between two successive major windows $u-1$ and u , both must use the same w ; otherwise, the comparison is not meaningful because the terrain may be totally different. Hence, after adapting w , we should wait at least two major windows before changing it again.

To understand how weights should be updated in Step 10, we examine all the possible behaviors of the search trajectory in successive major windows. We have identified four possible cases.

First, the trajectory does not stay within a feasible region, but goes from one feasible region to another through an infeasible region. During this time, $v_{\max}(i)$ is zero when the trajectory is in the first feasible region, increased when it travels from the first feasible region to an infeasible region, and decreased when going from the infeasible region to the second feasible region. No oscillations will be observed because oscillations normally occur around an equilibrium point in one feasible region. In this case, w is not changed.

Second, the trajectory oscillates around an equilibrium point of a feasible region. This can be detected when the number of oscillations in each major window is larger than a certain threshold, the trajectory is not always in a feasible region, and the trend of the maximum violation does not decrease. To determine whether the oscillations will subside eventually, we

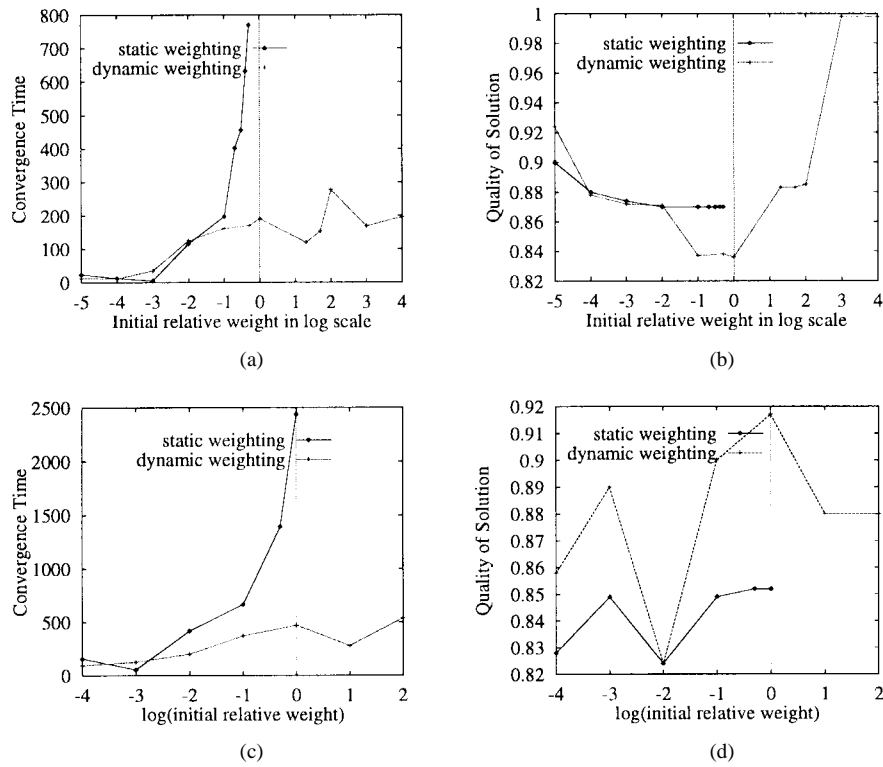


Fig. 6. Comparison of convergence time and quality of solution between static weighting and dynamic weighting for multiplierless QMF-bank design problems 32d and 48e, where quality is measured by the ratio of the reconstruction error of our design to that of Johnston’s design [10]. Hence, better designs have smaller values of quality. (a) Problem 32d with static weights. (b) Problem 32d with dynamic weights. (c) Problem 48e with static weights. (d) Problem 48e with dynamic weights.

compute $\bar{v}_u - \bar{v}_{u+1}$, the difference of the average values of maximum violation $v_{\max}(i)$ for two successive major windows u and $u + 1$. If the difference is not reduced reasonably, then we assume that the trajectory has not converged and decrease w accordingly.

Third, the search trajectory moves very slowly within a feasible region. This happens when w is very small, and the constraints dominate the search process. As a result, the objective value is improving very slowly and may eventually converge to a poor value. This situation can be identified when the trajectory remains within a feasible region in two successive major windows and is improving in successive major windows, but the improvement of the objective is not fast enough and is below an upper bound. Obviously, we need to increase w in order to speed up the improvement of the objective. If the objective remains unchanged, then the trajectory has converged, and no further modification of w is necessary.

Finally, the trajectory does not oscillate when it starts within a feasible region, goes outside the region, and converges to a point on the boundary. Here, a large w makes it more difficult to satisfy the constraints, causing the trajectory to move slowly to the feasible region. In this case, an appropriate decrease of w will greatly shorten the convergence time.

Table II illustrates the improvements in convergence times using adaptive weights. For all the initial weights considered, the adaptive algorithm is able to find converged designs in a reasonable amount of time, although the solution quality is not always consistent.

V. EXPERIMENTAL RESULTS

We have applied DLM-98 to solve the QMF-bank design problems formulated by Johnston [10]. In this section, we compare the performance of designs found by DLM-98 and those by Johnston [10], Chen *et al.* [4], Novel [27], simulated annealing (SA), and genetic algorithms (GA). All the experiments were run on Pentium Pro 200 computers with Linux unless specified otherwise.⁴

Our goal is to find designs that are better than the baseline results across all six performance measures. Hence, we use (5) with the constraint bounds defined by those of the baseline designs.

A. Performance of Lagrangian Methods with Dynamic Weights

To design multiplierless QMF banks, we allow the maximum number of ONE bits to be six and the minimum exponent to be -22 for each filter coefficient. The Lagrangian method uses both static weights and dynamic weights to solve 32d and 48e problems. We compare both the convergence time and the quality of solution in terms of reconstruction error. The starting points were obtained from Johnston’s design, and the control parameters were the same as those used in the previous subsection except that the window size N_u is 10.

A comparison of DLM-98 with static weights and that with dynamic weights is shown in Fig. 6. Even though the initial weights have very large ranges, $[10^{-5}, 10^4]$ for 32d and

⁴Filter-bank coefficients [Online]. Available FTP: manip.crhc.uiuc.edu/pub/papers/PostScript/J66/J66.coefficients

TABLE III
EXPERIMENTAL RESULTS OF DLM-98 IN SOLVING MULTIPLIERLESS QMF-BANK
DESIGN PROBLEMS. THE INITIAL POINTS OF THE RUN WERE FROM SIX
ONE-BIT EXPRESSIONS OF SCALED JOHNSTON'S SOLUTIONS

Filter	E_r	δ_p	E_p	δ_s	E_s	T_r	Scaling Factor	Time (hrs)
16a	0.99	0.99	0.94	0.99	0.95	0.99	0.9747	1.6
16b	0.99	0.99	0.90	0.99	0.98	0.99	0.8524	2.1
16c	0.96	0.99	0.98	0.99	0.99	0.99	0.5967	3.0
24b	0.97	0.99	0.87	0.96	0.99	0.99	0.9661	5.6
24c	0.89	0.99	0.58	0.99	0.99	0.99	0.6413	12.0
24d	0.81	0.99	0.83	0.99	0.99	0.99	0.5342	13.1
32c	0.96	0.99	0.75	0.99	0.99	0.99	0.5706	12.0
32d	0.83	0.95	0.61	0.99	0.99	0.99	0.6971	3.1
32e	0.72	0.99	0.90	0.99	0.99	0.99	0.5019	7.2
48c	0.88	0.95	0.85	0.99	0.99	0.99	0.7914	18.0
48d	0.95	0.99	0.75	0.99	0.99	0.99	0.7138	23.0
48e	0.91	0.99	0.80	0.99	0.99	0.99	0.5793	8.0
64d	0.87	0.99	0.83	0.76	0.99	0.99	0.9955	9.5
64e	0.85	0.99	0.73	0.99	0.99	0.99	0.8026	24.1

$[10^{-4}, 10^2]$ for 48e, the dynamic weight-adaptation algorithm converges in less than 300 min for the 32d problem and 510 min for 48e. However, using DLM-98 with static weights, when the initial w is larger than 1.0, the search cannot converge within 15 h for 32d and 32 h for 48e.

Note that, for 48e, the solution quality of DLM-98 with static weights is slightly better than our dynamic weight-adaptation algorithm for some initial weights. This happens because the latter may change the terrain during the search and find different solutions.

Finally, Table III shows the results of solving all the Johnston's benchmarks using filter coefficients with a maximum of six ONE bits. Our results show that we were able to find designs that have better reconstruction errors, while the other performance metrics are either the same or better.

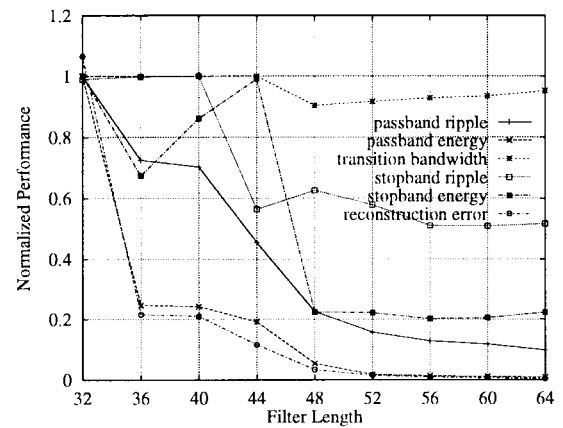
B. Comparison of DLM-98 with Johnston's Designs

In this section, we compare the performance of designs found by DLM-98 and those by Johnston [10].

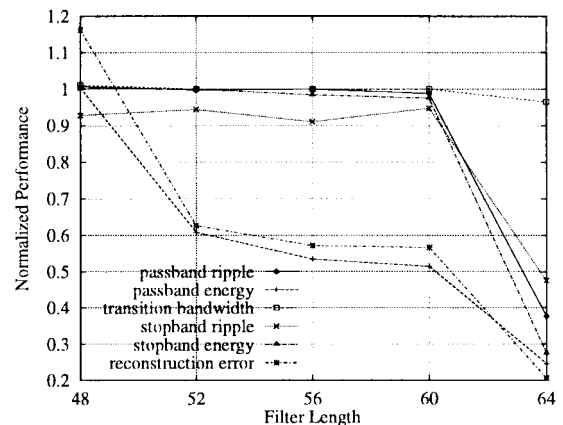
There are two parameters in a PO2 filter bank design: the maximum number of ONE bits in each filter coefficient and the number of filter taps. In our experiments, we have varied one while keeping the other fixed when evaluating a PO2 design with respect to a benchmark design.

We have used closed-form integration to compute the performance values. In contrast, Johnston [10] used sampling to compute energies. Hence, designs found by Johnston are not necessarily at the local minima in a continuous sense. To demonstrate this, we applied local search in a continuous formulation of the 24D design, starting from Johnston's design. We found a design with a reconstruction error of 3.83E-05, which is better than Johnston's result of 4.86E-05. By applying global search, we can further improve the design to have a reconstruction error of 3.66E-05.

We have evaluated PO2 designs obtained by DLM-98 with respect to Johnston's designs whose coefficients are 32-bit real numbers. Using the performance of Johnston's 32e design as constraints [10], we ran DLM-98 from ten different starting points obtained by randomly perturbing 1% of all the coefficients of Johnston's design [10]. Each run was limited so that each ONE bit of the coefficient was processed in



(a)



(b)

Fig. 7. Normalized performance for PO2 filter banks with a maximum of three ONE bits per coefficient and different number of filter taps. (a) Problem 32e. (b) Problem 48e.

a round-robin fashion 400 times. We then picked the best solution of the ten runs and plotted the result in Fig. 7, which shows the normalized performance of PO2 designs with increasing number of filter taps, while each filter coefficient has a maximum of three ONE bits. (The best design is one with the minimum reconstruction error if all the constraints are satisfied; otherwise, the one with the minimum violation is picked.) Our results show a design with 32 taps that is nearly as good as Johnston 32e's design. For filters with 32, 36, 40, and 44 taps, we used a starting point derived from Johnston's 32e design with filter coefficients first scaled by 0.5565 and truncated to a maximum of three ONE bits, and the filter coefficients of the remaining taps set to zeros initially. Starting points for filters with longer than 44 taps were generated similarly, except that a scaling factor of 0.5584 was used instead. Our results show that, as the filter length is increased, all the performance metrics improve, except the transition bandwidth, which remains close to that of the benchmark design.

With respect to Johnston's 48e design [10], we set a limit so that each ONE bit of the coefficient was processed in a round-robin fashion 800 times, and ran DLM-98 once from the truncated Johnston's 48e design. (The scaling factor was

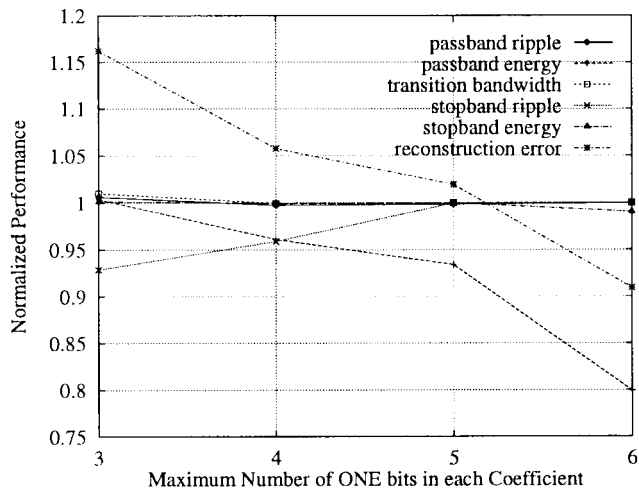


Fig. 8. Normalized performance with respect to Johnston's 48e QMF bank [10] for PO2 filters with 48 taps and different maximum number of ONE bits per coefficient.

0.5584 for filters with 48, 52, 56, and 60 taps. The scaling factor was 0.6486 for filters with 64 taps.) Our results show that our 48-tap PO2 design is slightly worse than that of Johnston's, while PO2 designs with 52 taps or longer have performance that are either the same or better than those of Johnston's 48e design. In particular, the reconstruction error of our 52-tap PO2 design is 62% of Johnston's 48e design, while that of our 64-tap PO2 design is only 21% of Johnston's 48e design.

In the next set of experiments, we kept the same number of taps as Johnston's 48e design and increased the maximum number of ONE bits in each coefficient from three to six. We set a limit so that each ONE bit of the coefficient was processed in a round-robin fashion 800 times, and ran DLM-98 once from the truncated Johnston's 48e design. Fig. 8 shows a design that is better than Johnston's 48e design when the maximum number of ONE bits per coefficient is six. In this case, the reconstruction error is 91% of Johnston's 48e design. (The scaling factors used are 0.5584 for three bits, 0.8092 for four bits, 0.7409 for five bits, and 1.0 for six bits.)

With respect to Johnston's 64d and 64e designs, Table IV shows improved PO2 designs obtained by DLM-98 using a maximum of six ONE bits per coefficient and 64 taps. No improvements were found when the maximum number of ONE bits is less than six.

C. Comparison of DLM-98 with Other Optimization Methods

In this section, we compare the performance of designs found by DLM-98 and those by Chen *et al.* [4], *Novel* [27], SA, and GA. Table IV shows improved designs found by DLM-98 with respect to Chen *et al.*'s designs with, respectively, 64 and 80 taps, all using a maximum of three ONE bits per coefficient. In these designs, we used Chen *et al.*'s designs as starting points and ran DLM-98 once with a limit so that each ONE bit was processed in a round-robin fashion 1000 times.

We also compare in Table IV the performance of 32e PO2 filter banks obtained by DLM-98 with a maximum of three ONE bits per coefficient, and those obtained by

TABLE IV

COMPARISON OF NORMALIZED PERFORMANCE OF FILTER BANKS WITH DISCRETE COEFFICIENTS DESIGNATED BY DLM-98 WITH RESPECT TO THOSE WITH CONTINUOUS COEFFICIENTS DESIGNATED BY JOHNSTON, CHEN, *NOVEL*, SIMULATED ANNEALING (SIMANN), AND GENETIC ALGORITHMS (EA-Ct AND EA-Wt).

Type	Discrete Coefficients					Continuous Coefficients			
	DLM-98					<i>Novel</i>	SA	EA-Ct	EA-Wt
Method	J-32e	J-64d	J-64e	C-64	C-80	J-32e	J-32e	J-32e	J-32e
Problem	J-32e	J-64d	J-64e	C-64	C-80	J-32e	J-32e	J-32e	J-32e
E_r	0.83	0.90	0.89	0.91	0.95	0.712	0.500	0.724	0.507
E_p	1.00	0.82	0.83	0.80	0.96	0.896	0.582	0.905	0.590
E_s	1.00	1.00	1.00	1.00	0.86	1.000	1.000	1.000	0.999
δ_p	1.00	0.97	1.00	1.00	1.00	1.000	1.000	1.000	0.997
δ_s	0.99	0.75	1.00	1.00	1.00	1.000	1.000	1.000	0.999
T_t	1.00	1.00	1.00	1.00	1.00	1.000	1.013	1.000	1.013

(Columns 2–4 show the performance of dlm-98 using three ONE bits for 32-tap filters and six ONE bits for 64-tap filters normalized with respect to that of Johnston's 32e, 64d, and 64e filter banks [10]. Columns 5–6 show the performance of dlm-98 using three ONE bits normalized with respect to that of Chen *et al.*'s 64-tap and 80-tap filter banks [4]. Columns 7–10 show the performance of 32-tap filter bank and using johnston's design as constraints.)

Novel, simulated annealing (SA), and evolutionary algorithms (EA's). *Novel* uses a continuous trace function to bring a search out of local minima rather than restarting the search from a new starting point when the search finds a feasible design. The SA we have used is SIMANN from netlib that works on a weighted-sum formulation. The EA is Sprave's Lice (linear cellular evolution) that can be applied to both constrained and weighted-sum formulations. SIMANN and EA-Wt use weighted-sum formulations with weight 1 for the reconstruction error and weight 10 for the remaining metrics. EA-Ct works on the same constrained formulation defined in (4). All methods were run significantly long with over 10 hours on a SUN SS20 workstation in each run.

We have tried various parameter settings and report the best solutions in Table IV. *Novel* improves Johnston's designs consistently. SIMANN and EA-Wt have difficulty in improving over Johnston's design across all measures and have found designs with larger transition bandwidth. EA-Ct found a design that improves Johnston's across all measures, although it is not as good as the one found by *Novel*. Note that all these designs have continuous coefficients that will need either a complex carry-save adder or a 32-bit multiplier in each tap in their hardware implementations. In contrast, DLM-98 obtained a design that improves E_r , while the other metrics are either exactly the same or slightly better than those of Johnston's. Moreover, the design uses a maximum of five additions in each tap, leading to very cost-effective implementations.

Since existing optimization packages like SIMANN and EA works in continuous space, we have also constructed our own simulated annealing package called *discrete simulated annealing* (DSA) that works directly in discrete space. As SA cannot handle constraints directly, we create a single objective based on a weighted sum of the objective and the constraints using static weights

$$F = w_0 V_{E_r} + w_1 V_{\delta_p} + w_2 V_{E_p} + w_3 V_{T_t} + w_4 V_{\delta_s} + w_5 V_{E_s}. \quad (29)$$

DSA first defines an initial temperature T_0 , and selects a starting point and scaling factor in the same way as that in Section IV-A. It then generates a new x' in discrete space and accepts the new point at the current temperature T according

TABLE V
EXPERIMENTAL RESULTS OF DSA IN DESIGNING MULTIPLIERLESS
QMF-BANK PROBLEM 24c, STARTING FROM A SIX ONE-BIT
EXPRESSION OF SCALED JOHNSTON'S SOLUTIONS

T_0	E_r	δ_p	E_p	δ_s	E_s	T_r	Search Time (Hours)
5.0	-	-	-	-	-	-	1.5
1.0	0.81	0.81	0.65	1.04	0.83	0.94	0.9
0.5	0.99	1.008	0.99	0.99	0.99	0.99	1.0
0.1	0.91	0.76	0.72	1.01	0.94	1.01	1.5
0.05	0.96	0.97	0.88	1.009	0.90	1.01	1.6
0.01	0.88	0.91	0.75	1.006	0.97	0.99	1.2
0.005	0.98	0.96	0.93	0.99	0.99	0.99	2.4
0.001	0.87	0.94	0.70	1.003	0.96	0.99	1.7

to the following probability:

$$\text{probability of accepting } x' = e^{-[(L(x')-L(x))^+/T]}$$

$$\text{where } a^+ = \begin{cases} a, & \text{if } a > 0 \\ 0, & \text{otherwise.} \end{cases} \quad (30)$$

Periodically, T is scaled down by scale_T when the maximum violation does not decrease over a period of time (set to ten round-robins in our experiments). Finally, DSA reports the best solution when the search converges.

In our experiments using DSA, we found it very difficult to set T_0 , scale_T , and the static weights in (29) that lead to better PO2 designs. A set of improperly chosen parameters will lead to violations of certain constraints. This phenomenon is obvious because the weights define the relative importance of the constraints.

Our experience on DSA is illustrated in the search of a better design of Johnston's 24c filter bank. After extensive experimentation, we initialized the weights to be $w_0 = 1.0$, $w_1 = w_2 = w_4 = 5.0$, $w_3 = 15.0$, $w_5 = 25.0$, and $\text{scale}_T = 0.95$. We further set the scaling factor to be 0.6413, the same as that in DLM-98 for 24c. Table V lists the eight designs found by DSA. When the initial temperature was too high (≥ 5.0), DSA did not find any meaningful design, but found near feasible designs when the initial temperature is lower. When the initial temperature is 0.005, DSA found a feasible PO2 design with six ONE bits that is slightly better than Johnston's 24c. Note that we did not find any feasible design after trying many other combinations of parameters.

In short, we found it difficult to use global search strategies, like SA and GA, to design PO2 filter banks formulated as weighted sum of the objective and the constraints. Without dynamically changing the weights as in DLM-98, it is hard to choose a proper set of weights (except by trial-and-error) that will allow SA or GA to converge to feasible designs. The best that SA and GA can find are designs with tradeoffs on different metrics. For this reason, the method studied in this paper represents a significant advance in solving discrete constrained optimization problems.

VI. CONCLUSION

We have presented a new discrete Lagrangian method (DLM-98) for designing multiplierless PO2 QMF banks. Our results show that DLM-98 can find better PO2 filter banks with very few ONE bits in each filter coefficient than other discrete

and continuous optimization methods. Our design method is unique because it starts from a constrained formulation, with the objective of finding a design that improves over a benchmark design. In contrast, existing methods for designing PO2 filter banks can only obtain designs with different tradeoffs among the performance metrics and cannot guarantee that the final design is always better than the benchmark design with respect to all the performance metrics.

REFERENCES

- [1] A. D. Booth, "A signed binary multiplication technique," *Quart. J. Mech. Appl. Math.*, vol. 4, pp. 236–240, 1951.
- [2] R. A. Caruana and B. J. Coffey, "Searching for optimal FIR multiplierless digital filters with simulated annealing," Philips Laboratories, Tech. Rep., 1988.
- [3] C.-K. Chen and J.-H. Lee, "Design of quadrature mirror filters with linear phase in the frequency domain," *IEEE Trans. Circuits Syst. II*, vol. 39, pp. 593–605, Sept. 1992.
- [4] ———, "Design of linear-phase quadrature mirror filters with powers-of-two coefficients," *IEEE Trans. Circuits Syst. I*, vol. 41, pp. 445–456, July 1994.
- [5] C. D. Creusere and S. K. Mitra, "A simple method for designing high-quality prototype filters for m-band pseudo QMF banks," *IEEE Trans. Signal Processing*, vol. 43, pp. 1005–1007, Apr. 1995.
- [6] M. H. Er and C. K. Siew, "Design of FIR filters using quadrature programming approach," *IEEE Trans. Circuits Syst. II*, vol. 42, pp. 217–220, Mar. 1995.
- [7] A. M. Geoffrion, "Lagrangian relaxation and its uses in integer programming," *Math. Prog. Study*, vol. 2, pp. 82–114, 1974.
- [8] B. R. Horng and A. N. Willon, Jr., "Lagrange multiplier approaches to the design of two-channel perfect-reconstruction linear-phase FIR filter banks," *IEEE Trans. Signal Processing*, vol. 40, pp. 364–374, Feb. 1992.
- [9] V. K. Jain and R. E. Crochiere, "Quadrature mirror filter design in the time domain," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-32, pp. 353–361, Apr. 1984.
- [10] J. D. Johnston, "A filter family designed for use in quadrature mirror filter banks," in *Proc. IEEE Proc. Int. Conf. ASSP*, 1980, pp. 291–294.
- [11] D. Kodek and K. Steiglitz, "Comparison of optimal and local search methods for designing finite word length FIR digital filters," *IEEE Trans. Circuits Syst.*, vol. CAS-28, pp. 28–32, Jan. 1981.
- [12] R. D. Koilpillai and P. P. Vaidyanathan, "A spectral factorization approach to pseudo-QMF design," *IEEE Trans. Signal Processing*, vol. 41, pp. 82–92, Jan. 1993.
- [13] Y. C. Lim and S. R. Parker, "FIR filter design over a discrete power-of-two coefficient space," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-31, pp. 583–519, June 1983.
- [14] D. G. Luenberger, *Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley, 1984.
- [15] K. Nayebi, T. P. Barnwell, III, and M. J. T. Smith, "Time-domain filter bank analysis: A new design theory," *IEEE Trans. Signal Processing*, vol. 40, pp. 1412–1429, June 1992.
- [16] T. Q. Nguyen, "Digital filter bank design quadratic-constrained formulation," *IEEE Trans. Signal Processing*, vol. 43, pp. 2103–2108, Sept. 1995.
- [17] H. Samuelli, "An improved search algorithm for the design of multiplierless FIR filters with powers-of-two coefficients," *IEEE Trans. Circuits Syst.*, vol. 36, pp. 1044–1047, July 1989.
- [18] J. D. Schaffer and L. J. Eshelman, "Designing multiplierless digital filters using genetic algorithms," in *Proc. Int. Conf. Genetic Algorithms*. San Mateo, CA: Morgan Kaufmann, 1993, pp. 439–444.
- [19] Y. Shang and B. W. Wah, "A discrete Lagrangian based global search method for solving satisfiability problems," *J. Global Optimization*, vol. 12, no. 1, pp. 61–99, Jan. 1998.
- [20] J.-J. Shyu and Y.-C. Lin, "A new approach to the design of discrete coefficient FIR digital filters," *IEEE Trans. Signal Processing*, vol. 43, pp. 310–314, Jan. 1995.
- [21] M. J. T. Smith and T. P. Barnwell, III, "Exact reconstruction techniques for tree-structured subband coders," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-34, pp. 434–441, June 1986.
- [22] I. Sodagar, K. Naybei, and T. P. Barnwell, III, "Time-varying filter banks and wavelets," *IEEE Trans. Signal Processing*, vol. 42, pp. 2983–2996, Nov. 1994.
- [23] A. K. Soman, P. P. Vaidyanathan, and T. Q. Nguyen, "Linear phase paraunitary filter banks: Theory, factorizations and designs," *IEEE Trans. Signal Processing*, vol. 41, pp. 3480–3496, Dec. 1993.

- [24] T. E. Tuncer and T. Q. Nguyen, "General analysis of two-band QMF banks," *IEEE Trans. Signal Processing*, vol. 43, pp. 544–548, Feb. 1995.
- [25] P. P. Vaidyanathan, *Multirate Systems and Filter Banks*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [26] B. W. Wah and Y. Shang, "A discrete Lagrangian-based global-search method for solving satisfiability problems," in *Satisfiability Problem: Theory and Applications*, D.-Z. Du, J. Gu, and P. Pardalos, Eds., Singapore American Mathematical Society, DIMACS Series in Discrete Mathematics and Theoretical Computer Science, 1997, pp. 365–392.
- [27] B. W. Wah, Y. Shang, T. Wang, and T. Yu, "Global optimization of QMF filter-bank design using NOVEL," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, Munich, Germany, Apr. 1997, vol. 3, pp. 2081–2084.
- [28] B. W. Wah, Y. Shang, and Z. Wu, "Discrete Lagrangian method for optimizing the design of multiplierless QMF filter banks, in *Proc. IEEE Int. Conf. Application-Specific Array Processors*, July 1997, pp. 529–538.
- [29] Z. Wu, "Discrete Lagrangian methods for solving nonlinear discrete constrained optimization problems," M.Sc. thesis, Dept. Comput. Sci., Univ. of Illinois, Urbana, IL, May 1998.



Benjamin W. Wah (S'74–M'77–SM'85–F'91) received the Ph.D. degree in computer science from the University of California at Berkeley in 1979.

He is currently the Robert T. Chien Professor of Engineering, a Professor in the Department of Electrical and Computer Engineering, Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, and a Professor at the Beckman Institute, University of Illinois at Urbana-Champaign. Previously, he had served on the faculty of Purdue University (1979–1985), as a Program

Director at the National Science Foundation (1988–1989), as Fujitsu Visiting Chair Professor of Intelligence Engineering, University of Tokyo (1992), and McKay Visiting Professor of Electrical Engineering and Computer Science, University of California at Berkeley (1994). His current research interests are in the areas of nonlinear search and optimization, knowledge engineering, multimedia signal processing, and parallel and distributed processing.

Dr. Wah was the Editor-in-Chief of the IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING between 1993–1996, and is the Honorary Editor-in-Chief of *Knowledge and Information Systems*. He currently serves on the editorial boards of *Information Sciences*, *International Journal on Artificial Intelligence Tools*, *Journal of VLSI Signal Processing*, and *Parallel Algorithms and Applications*. He had chaired a number of international conferences and is currently serving as the International Program Committee Chair of the IFIP World Congress in 2000. He had served in the IEEE Computer Society in various capacities and is currently the elected First Vice President for Publications. He is a Fellow of the Society for Design and Process Science. In 1989, he was awarded University Scholar of the University of Illinois, and in 1998, he received the IEEE Computer Society Technical Achievement Award.



Yi Shang (M'98) received the B.S. degree from the University of Science and Technology of China, the M.S. degree from the Chinese Academy of Sciences, Beijing, China, and the Ph.D. degree in computer science from the University of Illinois at Urbana-Champaign in 1997.

He joined the CECS department at the University of Missouri at Columbia as an Assistant Professor. His research interests are in the areas of artificial and computational intelligence, combinatorial and nonlinear optimization, and distributed and parallel computation. He has published many research papers in international journals and conferences. He has served on program committees for several IEEE and SPIE conferences and workshops.



Zhe Wu was born in Anhui, China, on January 4, 1974. He received B.Sc. degree in June 1996 from the Department of Special Class for Gifted Young, University of Science and Technology of China, and the M.S. degree in May 1998 from the Department of Computer Science, University of Illinois at Urbana-Champaign. He is currently working toward the Ph.D. degree at the University of Illinois at Urbana-Champaign.

His main research interests are in optimization, signal processing, software engineering and computer networks.