

# A Lattice Basis Reduction Approach for the Design of Finite Wordlength FIR Filters

Nicolas Brisebarre, Silviu-Ioan Filip, Guillaume Hanrot

► **To cite this version:**

Nicolas Brisebarre, Silviu-Ioan Filip, Guillaume Hanrot. A Lattice Basis Reduction Approach for the Design of Finite Wordlength FIR Filters. submitted for publication. 2016. <hal-01308801v2>

**HAL Id: hal-01308801**

**<https://hal.inria.fr/hal-01308801v2>**

Submitted on 30 May 2017 (v2), last revised 22 Feb 2018 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Lattice Basis Reduction Approach for the Design of Finite Wordlength FIR Filters

Nicolas Brisebarre, Silviu-Ioan Filip and Guillaume Hanrot

**Abstract**—Many applications of finite impulse response (FIR) digital filters impose strict format constraints for the filter coefficients. Such requirements increase the complexity of determining optimal designs for the problem at hand. We introduce a fast and efficient method, based on the computation of good nodes for polynomial interpolation and Euclidean lattice basis reduction. Experiments show that it returns quasi-optimal finite wordlength FIR filters; compared to previous approaches it also scales remarkably well (length 125 filters are treated in  $< 9$ s). It also proves useful for accelerating the determination of optimal finite wordlength FIR filters.

**Index Terms**—BKZ algorithm, Euclidean Lattice, finite wordlength, FIR coefficient quantization, FIR digital filters, HKZ algorithm, lattice basis reduction, LLL algorithm, minimax approximation

## I. INTRODUCTION

THE efficient design of FIR filters has been an active research topic in Digital Signal Processing (DSP) for several decades now. An important class of practical problems is related to the discrete nature of the representation space used for storing the filter coefficients (usually fixed-point or floating-point formats).

Typical constraints are that these coefficients be representable as signed  $b$ -bit integer values (up to a scaling factor) or that they have low complexity, such as them being sums of only a small number of power-of-two terms. While the first requirement is useful for implementation on fixed-point DSP processors, the second one occurs when one wants to construct multiplierless filtering circuits. It is addressed, for instance, in [1]–[9]. Our focus will be on the first requirement, in the case of direct-form FIR filters. This coefficient quantization issue has been considered as early as [10] (see also [11], the introduction of which gives an account on related work from the 1970s).

To better illustrate the problem, consider the following simple toy example. While not necessarily practical, it allows the reader to quickly grasp the context of our study. We want to compute a length 5 linear phase lowpass FIR filter with 5-bit wide fixed-point coefficients and uniformly weighted passband  $[0, 0.4\pi]$  and stopband  $[0.5\pi, \pi]$ . If we let  $D(\omega) = 1$  for  $\omega \in [0, 0.4\pi]$  and  $D(\omega) = 0$  if  $\omega \in [0.5\pi, \pi]$ , we want

$a_0, a_1$  and  $a_2 \in \{-2^4, \dots, 2^4 - 1\}$  such that the approximation error

$$\max_{\omega \in [0, 0.4\pi] \cup [0.5\pi, \pi]} \left| \frac{a_0}{2^4} + \frac{a_1}{2^3} \cos(\omega) + \frac{a_2}{2^3} \cos(2\omega) - D(\omega) \right|$$

is as small as possible.

The Parks-McClellan algorithm [12], for instance, outputs an optimal infinite precision frequency response  $H^*(\omega) = 0.430 \dots + 0.860 \dots \cos(\omega) + 0.069 \dots \cos(2\omega)$  and corresponding approximation error  $E^* = 0.360 \dots$ . If we round each coefficient towards the nearest number of the form  $k/2^3$  (or  $k/2^4$  for the constant coefficient), we obtain  $H_{\text{naive}}(\omega) = \frac{7}{2^4} + \frac{7}{2^3} \cos(\omega) + \frac{1}{2^3} \cos(2\omega)$ , and a corresponding error  $E_{\text{naive}} = 0.4375$ . On the other hand, the approach we introduce in the subsequent sections returns  $H(\omega) = \frac{8}{2^4} + \frac{6}{2^3} \cos(\omega) + \frac{1}{2^3} \cos(2\omega)$ , which improves the error to  $E = 0.375$ . Using a mixed integer linear programming (MILP) routine, we obtain that  $H$  is actually an optimal 5-bit coefficient frequency response  $H_{\text{opt}}$  and  $E$  is the optimal approximation error  $E_{\text{opt}}$ .

This example displays two interesting facts:

- The straightforward approach, which we call “naive rounding”, consisting in rounding each coefficient of the infinite precision response to the nearest coefficient fitting the imposed constraints, may, even in very simple cases, yield results which are far from optimal (for instance, [13] mentions an example where a 30 dB improvement is observed).
- Another natural approach is to consider all the possible responses that we get by rounding up or down the coefficients of the infinite precision response. A first major drawback of this idea is that the number of possibilities is exponential in the degree of the filter. Moreover, our toy example shows that it is possible that none of them yields an optimal response, as was already noticed in [4].

In general, finding an optimal finite wordlength direct form FIR filter proves to be computationally expensive as the degree increases, with exact solvers making use of MILP techniques, sometimes combined with clever branch-and-bound strategies [1]–[3], [5], [11], [13]–[17]. To successfully cope with larger degrees, faster heuristics have also been proposed. They produce results that are, in many cases, quasi-optimal and can help speed up exact solvers. For instance, the approach introduced in [18], based on error spectrum shaping techniques, is routinely used inside MATLAB<sup>TM</sup> for FIR coefficient quantization. The methods developed in [4] (see also [19] for a variant) and [20] also produce very good results.

N. Brisebarre and G. Hanrot are with CNRS, ÉNS de Lyon, Inria, Université Claude Bernard Lyon 1, Laboratoire LIP (UMR 5668), 15 parvis René-Descartes, F-69007 Lyon, France.

E-mail: nicolas.brisebarre@ens-lyon.fr, guillaume.hanrot@ens-lyon.fr

S.-I. Filip is with CNRS, ÉNS de Lyon, Inria, Université Claude Bernard Lyon 1, Laboratoire LIP (UMR 5668), Lyon, France and Mathematical Institute, University of Oxford, Oxford, UK

E-mail: filips@maths.ox.ac.uk

This work was partly supported by the FastRelax project of the French Agence Nationale de la Recherche.

### A. Sketch of our approach and outline of the paper

The rest of the paper introduces a novel method for designing (quasi-)optimal direct form FIR filters with fixed-point and/or floating-point coefficients. It is based on previous work for doing machine-efficient polynomial approximation of functions [21] combined with techniques that yield families of very good interpolation nodes [22]–[27]. It turns out to be quite robust, especially when dealing with high degree designs. We also provide an open source C++ implementation of our approach<sup>1</sup>. Based on the output quality, we think that our method can also help accelerate exact solvers based on MILP.

As we will show in Section II, the problem that we address can be formulated as an MILP question, constructed from a sufficiently dense discretization of the target filter bands. The first key aspect of our work is the modeling of the MILP instance as a Closest Vector Problem (CVP), a fundamental question in the study of Euclidean lattices<sup>2</sup> [28]. The latter is a discrete structure which benefits from rich algorithmic results that we take advantage of. In particular, there exist very efficient algorithms for computing approximate (and yet, at least in the filter design setting, often close to optimal) solutions to a CVP.

The CVP instance that we use is also constructed from a discretization of the filter frequency bands. For it to be effective, two constraints are imposed:

- the number of discretization points should be as small as possible, both to produce a tractable CVP question and for a technical reason to be made clear in Section V;
- the discretization points should be chosen so that the CVP instance models as faithfully as possible the initial MILP problem.

To determine such a discretization, a second key idea is to compute a certain family of interpolation nodes, which gives rise to excellent polynomial approximations on the bands of the filter. These points are analogous to the so-called Chebyshev nodes on a closed interval.

Then, we use an idea due to R. Kannan [29] to compute a good approximate solution to this CVP. A simple additional trick from [21] applied to this solution yields a certain frequency response  $H$ : in practice, its coefficients satisfy the requested constraints and offer a very good approximation to the ideal frequency response.

After giving a precise statement of the problem that we want to solve in Section II, Section III presents the discretizations that we use, while Section IV recalls the algorithmic questions attached to Euclidean lattices that we face in our approach and some state-of-the-art algorithms to address them. They are later used for detailing our method in Section V. Examples and a comparison to other approaches are discussed in Section VI, followed by concluding remarks in Section VII.

## II. THE FINITE WORDLENGTH SETTING

We first recall some well-known ideas related to linear-phase FIR filter design [30]–[32]<sup>3</sup>. Such a process is typically carried

<sup>1</sup>See <https://github.com/sfilip/fquantizer>.

<sup>2</sup>Note that there is no link between Euclidean lattices and lattice wave digital filters, as the one used in [6] for instance.

<sup>3</sup>A part of this section follows Section II of [17].

out in the frequency domain. In case of an  $N$ -tap causal filter with real valued impulse response  $\{h[n]\}$ , the corresponding frequency response we want to optimize is

$$\begin{aligned} H(\omega) &= \sum_{k=0}^{N-1} h[k] e^{-i\omega k} \\ &= G(\omega) e^{i\left(\frac{L\pi}{2} - \frac{N-1}{2}\omega\right)}, \end{aligned}$$

where  $G$  is a real valued function and  $L$  is 0 or 1. Traditionally, there are four such types of FIR filters considered in practice, labeled from I through IV; they depend on the parity of the filter length  $N$  and on the symmetry of  $h$  (positive for  $L = 0$  and negative for  $L = 1$ ). We can express  $G(\omega)$  as  $Q(\omega)P(\omega)$ , where  $P(\omega)$  is of the form:

$$P(\omega) = \sum_{k=0}^n p_k \cos(k\omega).$$

Depending on the filter type,  $Q(\omega)$  is either 1,  $\cos(\omega/2)$ ,  $\sin(\omega)$  or  $\sin(\omega/2)$ . We have  $n = \lfloor N/2 \rfloor$  and there are simple formulas linking the  $h[k]$ 's and the  $p_k$ 's (see for instance [31, Ch. 15.8–15.10]).

**Remark 1.** *If we consider the change of variable  $x = \cos(\omega)$ ,  $P(\omega) = \sum_{k=0}^n p_k \cos(k\omega)$  is in fact a polynomial of degree at most  $n$  in  $\cos(\omega)$  expressed in the basis of Chebyshev polynomials  $(T_k)_{0 \leq k \leq n}$  of the first kind, i.e.,*

$$P(\omega) = \sum_{k=0}^n p_k T_k(\cos(\omega)) = \sum_{k=0}^n p_k T_k(x).$$

*It will sometimes be convenient to see our problems in the language of polynomial approximation. When this is the case,  $x$  will be the approximation variable and  $X \subseteq [-1, 1]$  the transformed domain associated to  $\Omega$  i.e.,  $\cos \Omega$ .*

While the focus of our presentation is on type I filters (hence  $Q \equiv 1$ ), adapting it to the other three cases is straightforward.

The optimal length  $N$  frequency response with infinite precision coefficients is the solution of the following:

**Problem 1** (Equiripple (or minimax) FIR filter design). *Let  $\Omega$  be a compact subset of  $[0, \pi]$  and  $D(\omega)$  an ideal frequency response, continuous on  $\Omega$ . For a given filter degree  $n \in \mathbb{N}$ , we want to determine  $P^*(\omega) = \sum_{k=0}^n p_k^* \cos(\omega k)$  such that the weighted error function  $E^*(\omega) = W(\omega) (P^*(\omega) - D(\omega))$  has **minimum** uniform norm*

$$\|E^*\|_{\infty, \Omega} = \sup_{\omega \in \Omega} |E^*(\omega)|,$$

*with the weight function  $W$  continuous and positive over  $\Omega$ .*

There exists a unique solution to this problem [33, Ch. 3.4]:

**Theorem 1** (Alternation theorem). *A necessary and sufficient condition for  $P(\omega)$  to be the **unique** transfer function of degree at most  $n$  that minimizes the weighted approximation error  $\delta_\Omega = \|E(\omega)\|_{\infty, \Omega}$  is that  $E(\omega) = W(\omega) (P(\omega) - D(\omega))$  has **at least**  $n + 2$  equioscillating extremal frequencies over  $\Omega$ ; i.e., there exist at least  $n + 2$  values  $\omega_k$  in  $\Omega$  such that  $\omega_0 < \omega_1 < \dots < \omega_{n+1}$  and*

$$E(\omega_k) = -E(\omega_{k+1}) = \lambda(-1)^k \delta_\Omega, \quad k = 0, \dots, n,$$

where  $\lambda \in \{\pm 1\}$  is fixed.

The usual way to solve Problem 1 is to use, as already mentioned in Section I, the Parks-McClellan version of the Remez algorithm.

The finite wordlength version of Problem 1 that is of interest here, and also in practice, proves to be much harder to solve in general. For fixed-point coefficient quantization we can basically consider the impulse response coefficients to be scaled integers. If we want to use  $b$ -bit fixed-point values, we view them under the form  $m/s$ , where  $m$  is an integer in  $I_b = \{-2^{b-1}, \dots, -1, 0, 1, \dots, 2^{b-1} - 1\}$  and  $s$  is a *fixed* scaling factor, which, in many cases, is a power of 2. Let  $\varphi_0(\omega) = 1/s$  and  $\varphi_k(\omega) = 2 \cos(k\omega)/s$  for  $k = 1, \dots, n$ . We then have:

**Problem 2** (Finite wordlength minimax approximation). *Consider  $\Omega, D, n$  and  $W$  defined as in Problem 1. Determine  $P(\omega) = \sum_{k=0}^n m_k \varphi_k(\omega)$ , where  $m_k \in I_b$  for  $k = 0, \dots, n$ , such that the error*

$$\sup_{\omega \in \Omega} |W(\omega) (P(\omega) - D(\omega))| \quad (1)$$

is *minimal*.

This problem always has a (not necessarily unique) solution. Expressed as an optimization question, it becomes:

$$\begin{aligned} & \text{minimize} && \delta \\ & \text{subject to} && W(\omega) \left( \sum_{k=0}^n m_k \varphi_k(\omega) - D(\omega) \right) \leq \delta, \omega \in \Omega, \\ & && W(\omega) \left( D(\omega) - \sum_{k=0}^n m_k \varphi_k(\omega) \right) \leq \delta, \omega \in \Omega, \\ & && m_k \in I_b, k = 0, \dots, n. \end{aligned} \quad (2)$$

In practice,  $\Omega$  is generally replaced with a finite discretization  $\Omega_d$ , giving rise to a MILP instance. A number of points equal to  $16n$  usually suffices, according to [34].

The scaling factor  $s$ , interpreted as the filter gain, can also be a design parameter. Optimizing for  $s$  as well can improve the quality of the quantization error (1) by a small factor. The caveat is that finding the best  $s$  is nontrivial [13], [34], [35].

More general quantization problems where coefficients have different formats from one another (be it fixed-point or floating-point) can also be expressed using the framework of Problem 2, with the difference being that each coefficient will have its own power of 2 scaling factor  $s_k, k = 0, \dots, n$ . Nevertheless, for simplicity, in the sequel we will only consider fixed-point quantization examples with a uniform scaling factor  $s$ .

### III. DISCRETIZATION OF THE PROBLEM

The first step of our approach is to discretize Problem 2 with (almost) as few points as possible, in order to speed up computation and also because our lattice-based approach requires a coarse discretization in order to return relevant results (we will make this point clear at the beginning of Section V). We pick  $\ell$  points  $\omega_0, \dots, \omega_{\ell-1}$  in  $\Omega$ , with  $\ell = n+1$  or slightly larger than  $n+1$ , and we want to determine a

$P(\omega) = \sum_{k=0}^n m_k \varphi_k(\omega)$ , where  $m_k \in I_b$  for  $k = 0, \dots, n$ , such that the error

$$\sup_{j=0, \dots, \ell-1} |W(\omega_j) (P(\omega_j) - D(\omega_j))|$$

is as small as possible. A question immediately arises: how can we force this discretized version to be as faithful as possible to the initial Problem 2?

Choosing appropriate points  $\omega_i$  (or  $x_i = \cos(\omega_i)$ ) is indeed critical for obtaining very good results. The preferred choices in [21] are the zeros of  $E^*(\omega)$  when first solving Problem 1 on a closed interval  $X = [a, b]$  or appropriately scaled  $n$ -th order Chebyshev nodes of the first kind

$$\mu_k = \frac{a+b}{2} + \frac{b-a}{2} \cos\left(\frac{(k+1/2)\pi}{n+1}\right), k = 0, \dots, n.$$

Such points are excellent when doing polynomial approximation (by interpolation) on a closed interval [36]. In the filter design setting however,  $\Omega$  (and consequently  $X = \cos \Omega$ ) is routinely a union of two or more closed intervals, and closed form expressions for appropriate  $\omega_i$ 's (or  $x_i$ 's) do not seem to be readily available. We now propose three complementary alternatives that work well together in our practical experimentations. The idea underlying the first two choices is to force the finite wordlength transfer function to mimic the minimax approximation  $P^*$ , whereas the third one corresponds to a relevant choice for the initialization of the Parks-McClellan algorithm to determine  $P^*$ .

#### A. Alternating extrema of the minimax error function

Inspired by the characterization provided by Theorem 1, the discretization  $\Omega_e$  we suggest is a set of  $n+2$  frequencies that yield  $n+2$  alternating extrema for the error function  $E^*$ .

There are some cases where the number of frequencies corresponding to these alternating extrema is larger than  $n+2$ . Therefore, we need an effective way to select exactly  $n+2$  among them: the Parks-McClellan algorithm iteratively constructs such a list.

**Remark 2.** *In practice, we use the implementation of the Parks-McClellan algorithm presented in [37] and available from <https://github.com/sfilip/firpm>: it computes  $P^*$  and a list of  $n+2$  nodes where  $E^*$  equalternates, our  $\Omega_e$ .*

#### B. “Zeros” of the minimax error function

This strategy is the analogue of one of the choices from [21]. And yet, there is a significant difference in our setting. Actually, in the case of a single interval  $[a, b]$ , Theorem 1 and the intermediate value theorem ensure that there exist, at least,  $n+1$  points  $x_0, \dots, x_n$  in  $[a, b]$  such that  $E^*(x_i) = 0$  for  $i = 0, \dots, n$ . Unfortunately, we do not have any such guarantee in the present multi-interval setting: the number of zeros might be less than  $n+1$  so we complete the list by picking points that are half-way the endpoints of consecutive bands. Note that the latter points do not belong to  $\Omega$  and we will explain in Section V how we use them. We execute the following:



**Algorithm 1** “Zeros” of the minimax error discretization

**Input:** the minimax approximation  $P^*$ , ideal response  $D$ , weight  $W$ , transformed domain  $X$

**Output:** an  $\ell \geq n+1$ -element discretization  $\Omega_z$  of  $[0, \pi]$  made of zeros of  $E^*$  and, if necessary, points that are half-way the extremities of the bands making up  $X = \cos \Omega$

// Take the zeros of  $E^*$  over  $X$

1:  $X_z \leftarrow \text{roots}(E^*, X)$

// if the number of zeros is less than  $n+1$ ,

// take the midpoints of the transition bands

// making up  $X$  (i.e.,  $X = \cup_{j=1}^k [x_0^{(j)}, x_1^{(j)}]$ )

// with  $x_1^{(j)} < x_0^{(j+1)}$  for  $j = 0, \dots, k-1$

2:  $j \leftarrow 1$

3: **while**  $|X_z| < n+1$  **do**

4:  $X_z \leftarrow X_z \cup \left\{ (x_1^{(j)} + x_0^{(j+1)})/2 \right\}$

5:  $j \leftarrow j+1$

6: **end while**

7:  $\Omega_z \leftarrow \arccos(X_z)$

**Remark 3.** One can prove that if a filter has  $p$  bands then  $E^*$  has at least  $n+2-p$  zeros and at most  $n+p-1$  zeros on the bands ( $p \geq 2$ ). Therefore, since there are  $p-1$  points that are half-way the extremities of the  $p$  bands, we are sure that the output discretization has at least  $n+1$  points and no more than  $n+p-1$  points. In the practical cases presented in Section VI, this discretization has exactly  $n+1$  points (2 band case) or at most  $n+2$  points (3 band case).

### C. Approximate Fekete points

The second author has previously used this last approach as a numerically robust way of initializing the Parks-McClellan algorithm [37, Sec. 4.1–4.2]. We now briefly recall it.

For any compact set  $X \subset \mathbb{R}$ , a set of Fekete points of degree  $n+2$  are the elements of a set  $\{t_0, \dots, t_{n+1}\} \subset X$  which maximize the absolute value of the Vandermonde-like determinant  $|w(y_i)T_j(y_i)|_{0 \leq i, j \leq n+1}$ , with  $w(x) = W(\arccos(x))$ , where  $W$  is the weight function used in the statements of Problems 1 and 2 and the  $T_j$ 's are the first  $n+2$  Chebyshev polynomials of the first kind.

Unfortunately, computing them is difficult in general. Still, *approximate* versions of such nodes, called Approximate Fekete Points (AFP), can be determined efficiently. The idea is to replace  $X$  by a suitable discretization  $Y \subseteq X$  with  $m+1 \geq n+2$  elements and extract Fekete points of  $Y$ . One good choice is that of so-called *weakly admissible meshes* [23]. The example we use consists of taking  $Y$  to be the union of the  $(n+1)$ -th order Chebyshev nodes of second kind

$$\nu_k = \frac{a+b}{2} + \frac{b-a}{2} \cos\left(\frac{k\pi}{n+1}\right), k = 0, \dots, n+1,$$

scaled to each interval that makes up  $X$ . Even in this case, obtaining Fekete points over  $Y$  is an NP-hard problem [38]. We can, however, use an effective greedy algorithm based on column pivoting QR [22], [24]–[26]:

**Algorithm 2** Approximate Fekete points discretization

**Input:** discretized subset  $Y = \{y_0, \dots, y_m\} \subseteq X$

**Output:** an  $n+2$  distinct element set  $X_{\text{afp}} \subset Y$  s.t. the Vandermonde-like submatrix  $\mathbf{V}(X_{\text{afp}})$  generated by the elements of  $X_{\text{afp}}$  has a large volume

// Initialization.  $\mathbf{V}(Y)$  is the Vandermonde-like matrix

//  $(w(y_i)T_j(y_i))_{0 \leq i \leq m, 0 \leq j \leq n+1}$

1:  $\mathbf{A} \leftarrow \mathbf{V}(Y)$ ,  $b \in \mathbb{R}^{n+2}$ ,  $b \leftarrow (1 \dots 1)^t$

// QR-based Linear System Solver

2: using a column pivoting-based QR solver (via an equivalent of LAPACK's DGEQP3 routine [39]), find  $\mathbf{w} \in \mathbb{R}^{m+1}$ , a solution to the underdetermined system  $\mathbf{b} = \mathbf{A}^t \mathbf{w}$ .

// Subset Selection

3: take  $X_{\text{afp}}$  as the set of elements from  $Y$  whose corresponding terms inside  $\mathbf{w}$  are different from zero, that is, if  $y_i \in Y$  and  $w_i \neq 0$ , then  $y_i \in X_{\text{afp}}$ .

Similar to the previous two choices, we will in fact work with  $\Omega_{\text{afp}} = \arccos(X_{\text{afp}})$ .

**Remark 4.** There are solid theoretical arguments for the AFP approach. They are based on the study of so-called Lebesgue constants and can be found in [37], [40].

**Remark 5.** Another alternative of good discretization grids is to take them according to the so-called equilibrium distribution for weighted minimax approximation on  $X$ . Articles like [41]–[43] touch on this aspect, with [42], [43] being particularly focused on the implications to FIR filter design problems. Computation generally involves the use of tools for numerical conformal mapping. We have not tested this idea here because it is more involved to set up than the approximate Fekete points method and it is not as flexible either, since it only deals with piecewise constant weight functions.

## IV. THE CLOSEST VECTOR PROBLEM AND LATTICE BASIS REDUCTION

The following four elements play a key role in this section and our approach. First, let  $\|\cdot\|_2$  and  $\|\cdot\|_\infty$  denote the usual Euclidean and supremum norms over  $\mathbb{R}^\ell$  (i.e.,  $\|\mathbf{v}\|_2 = \left(\sum_{i=0}^{\ell-1} v_i^2\right)^{1/2}$  and  $\|\mathbf{v}\|_\infty = \max_{i=0}^{\ell-1} |v_i|$ ).

The Euclidean lattice structure is a fundamental component of various problems in Mathematics, Computer Science or Crystallography [28], [44]–[48]:

**Definition 1.** A lattice  $L \subset \mathbb{R}^\ell$  is the set of integer linear combinations of a family  $(\mathbf{b}_1, \dots, \mathbf{b}_d)$  of  $\mathbb{R}$ -linearly independent vectors of  $\mathbb{R}^\ell$ . We shall then say that  $(\mathbf{b}_i)_{1 \leq i \leq d}$  is a **basis** of  $L$ , and that  $d(\leq \ell)$  is the **dimension** of  $L$ .

Finally, we define one of the most important algorithmic problems concerning Euclidean lattices:

**Problem 3** (Closest vector problem (CVP)). *Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^\ell$ . Given as input a basis of a lattice  $L \subset \mathbb{R}^\ell$  of dimension  $d$ ,  $d \leq \ell$ , and  $x \in \mathbb{R}^\ell$ , find  $y \in L$  s.t.  $\|x - y\| = \min\{\|x - v\|, v \in L\}$ .*

In our approach, in order to find a good solution to Problem 2, we address, after a suitable discretization, System (2). Up to

TABLE I  
LATTICE BASIS REDUCTION ALGORITHMS

Algorithm	First length defect	Time
LLL $_{\eta}$ [52]	$(\eta - 1/4)^{-(d-1)/2}$	$\text{poly}(d, \text{size}(B))$
BKZ $_{\beta}$ [53], [54]	$\approx \beta^{(d-1)/(\beta-1)}$	$\text{poly}(d, \text{size}(B))2^{O(\beta)}$
HKZ [55]	1	$2^d \text{poly}(d, \text{size}(B))$

a very mild relaxation of the constraint  $m_k \in I_b$  to  $m_k \in \mathbb{Z}$ , the latter is actually a CVP in the supremum (or  $\|\cdot\|_{\infty}$ ) norm. Although  $\text{CVP}_{\infty}$  can be attacked using MILP techniques, we propose to approximate it by the same CVP, but in the Euclidean (or  $\|\cdot\|_2$ ) norm, a problem which we shall denote by  $\text{CVP}_2$ . Even though, from a complexity-theoretic point of view, both problems are NP-hard [49], this change of perspective is motivated by the existence of a wealth of practical (i.e., efficient and with good performances) approximation algorithms to deal with this kind of problems.

The algorithms addressing a  $\text{CVP}_2$  usually rely on a preprocessing step called *lattice basis reduction*, a notion that we briefly review in Subsection IV-A, together with three main tools for performing it: the LLL, BKZ and HKZ algorithms.

As a strategy to solve  $\text{CVP}_2$ , we shall try to reduce it to finding short nonzero vectors in another lattice, a similar task where lattice basis reduction is crucial. Such reductions are well studied in the literature, see eg. [29], [50], [51]. They are usually based on a subtle use of the so-called Kannan embedding technique, which we also use, albeit in a simplified, but more efficient way (see Section IV-B). Experiments show this to suffice for the problem under study.

We also point out that the LLL algorithm has already been used in the digital filter design context, but in a different way [17]: in order to accelerate the MILP search for an optimal filter, Kodek determines, thanks to LLL<sup>4</sup>, a lower bound on the optimal approximation error, which is used as a first guess for  $\delta$  in (2). Moreover, to the best of our knowledge, the present text discusses the first use of the BKZ, HKZ and Kannan embedding algorithms for a filter design purpose.

#### A. Lattice basis reduction

A lattice of dimension  $d \geq 2$  has infinitely many bases. Among those bases, the ones which are made up of short and somewhat orthogonal vectors are usually considered good, whereas the other ones are deemed bad. Lattice basis reduction algorithms allow one to move from a bad basis to a good one, an important preconditioning step in solving lattice problems.

We shall consider three different lattice basis reduction algorithms, namely LLL, BKZ and HKZ. We shall not describe them nor their complete output – this is a rather technical topic and is out of the scope of this paper; we refer the reader to [28] for a detailed discussion of lattice basis reduction algorithms. Note however that the LLL algorithm depends on a real parameter  $\eta \in (1/4, 1)$ , while the BKZ algorithm depends on an integer parameter  $\beta \in [2, d]$ , the blocksize. For  $\beta = 2$ ,

BKZ is akin to LLL, whereas for  $\beta = d$ , BKZ is identical<sup>5</sup> to HKZ.

The performances of these algorithms are usually measured by the so-called first length defect, namely  $\|v_1\|_2/\lambda_1(L)$ , where  $v_1$  is the first vector of the basis output by the algorithm and  $\lambda_1(L)$  the length of a shortest non-zero vector.

As we shall use lattice basis reduction as a way to find short vectors, we sum up in Table I the performance, in terms of quality and speed, of LLL, BKZ and HKZ. Note that  $\text{size}(B)$  denotes the size of the input basis. The bottomline is that in terms of quality,  $\text{LLL} \leq \text{BKZ} \leq \text{HKZ}$ , while in terms of efficiency, unsurprisingly,  $\text{LLL} \geq \text{BKZ} \geq \text{HKZ}$ .

**Remark 6.** *It should be pointed that all these bounds are worst-case estimates, with little practical relevance in our specific context. Indeed, it seems that our input lattices and lattice bases are much better conditioned than expected, as all these algorithms seem to have a much better performance than expected on both quality and efficiency – HKZ-reducing a random lattice of dimension greater than 60 is, as of today, a considerable task, while we were able to do it on instances of degree 62 resulting from our problems in less than two minutes (see Section VI).*

#### B. The CVP problem & Kannan's embedding

When facing a  $\text{CVP}_2$  instance, one may choose to resort to exact, exponential-time approaches [56]. Even though they are efficient enough to deal with moderate size problems, we have found them to provide results which are not superior to our more efficient method. This might be explained by two reasons :

- as already observed, the lattice problems raised by digital filter design seem to be very well-conditioned, so that approximation algorithms tend to give much better results than expected (actually, often, optimal results).
- since  $\text{CVP}_2$  is already an approximate formulation of the problem to be solved, seeking an exact solution for it is not very meaningful. Actually, our main claim is that very good solutions to  $\text{CVP}_2$  correspond to very good solutions to (2) (and Section VI provides evidence for that), whereas experience shows that optimal  $\text{CVP}_2$  solutions do not, in general, correspond to optimal solutions of (2).

We use Kannan's embedding technique [29] in the following way: given a basis  $(\mathbf{b}_1, \dots, \mathbf{b}_d)$  of  $L$  and the target vector  $\mathbf{t}$ , we form the  $(d+1)$ -dimensional lattice  $L'$  of  $\mathbb{R}^{d+1}$  generated by the vectors  $\mathbf{b}'_i := \begin{pmatrix} \mathbf{b}_i \\ 0 \end{pmatrix}$  and  $\mathbf{t}' := \begin{pmatrix} \mathbf{t} \\ \gamma \end{pmatrix}$ , where  $\gamma$  is a real parameter to be chosen later on. If  $\mathbf{u}$  is a short vector of the lattice  $L'$ , there exist  $u_1, \dots, u_d, u_{d+1}$  such that  $\mathbf{u} = \sum_{k=1}^d u_k \mathbf{b}'_k + u_{d+1} \mathbf{t}'$ ; hence

$$\|\mathbf{u}\|_2^2 = \left\| \sum_{k=1}^d u_k \mathbf{b}_k + u_{d+1} \mathbf{t} \right\|_2^2 + \gamma^2 u_{d+1}^2.$$

<sup>5</sup>Note that the reason this statement seems to contradict Table I is that real BKZ estimates showing this behavior are technical and have been strongly simplified in Table I.

<sup>4</sup>Kodek actually applies LLL to a lattice which is dual to ours – which is coherent with the use he makes of it.

Since  $\mathbf{u}$  is assumed to be short (for instance, a vector in a reduced basis of  $L'$ ),  $\|\sum_{k=1}^d u_k \mathbf{b}_k + u_{d+1} \mathbf{t}\|_2^2$  is small. If  $u_{d+1} = \pm 1$ ,  $\mp \sum_{k=1}^d u_k \mathbf{b}_k \in L$  is close to the vector  $\mathbf{t}$ . In other words, if we are capable of finding a basis of short vectors of  $L'$  (the exact task that lattice basis reduction performs), we expect to obtain a vector in  $L$  very close to  $\mathbf{t}$ .

In our experiments, the (heuristic) choice  $\gamma = \max_{k=1}^d \|\mathbf{b}_k\|_2$  always yielded a vector with  $u_{d+1} = \pm 1$  in all cases (see [40, Ch. 4.3.4] for a theoretical justification). The entire procedure is summarized in the next listing:

---

### Algorithm 3 Kannan embedding

---

**Input:** basis  $(\mathbf{b}_1, \dots, \mathbf{b}_d)$  of  $L \subset \mathbb{R}^\ell$ , target vector  $\mathbf{t} \in \mathbb{R}^\ell$ , a lattice basis reduction algorithm (LLL, BKZ or HKZ)

**Output:**  $\mathbf{m} = (m_1 \dots m_d)^t \in \mathbb{Z}^d$  s.t.  $\sum_{k=1}^d m_k \mathbf{b}_k \approx \mathbf{t}$ , change of basis matrix  $\mathbf{U} \in \mathbb{Z}^{d \times d}$

```

1:  $\gamma \leftarrow \max_{1 \leq k \leq d} \|\mathbf{b}_k\|_2$ 
   // Construct the embedded basis
2:  $\mathbf{B}' \leftarrow \begin{pmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_d & \mathbf{t} \\ 0 & 0 & \dots & 0 & \gamma \end{pmatrix}$ 
   // Perform basis reduction (LLL, BKZ or HKZ) on  $\mathbf{B}'$ 
   //  $\mathbf{C}' = (c'_{i,j}) \in \mathbb{R}^{(\ell+1) \times (d+1)}$  is the reduced basis
   //  $\mathbf{U}' = (u'_{i,j}) \in \mathbb{Z}^{(d+1) \times (d+1)}$  is the change of basis
   // matrix (i.e.,  $\mathbf{C}' = \mathbf{B}'\mathbf{U}'$ )
3:  $(\mathbf{C}', \mathbf{U}') \leftarrow \text{LatticeReduce}(\mathbf{B}')$ 
   // Get change of basis matrix for  $\mathbf{B}$  (i.e., first  $d$ 
   // rows and columns of  $\mathbf{U}'$ )
4:  $\mathbf{U} \leftarrow (u'_{i,j})_{1 \leq i \leq d, 1 \leq j \leq d}$ 
   // Extract element from last line and column of  $\mathbf{C}'$ 
5:  $\gamma' \leftarrow c'_{\ell+1, d+1}$ 
   // Construct the outputs in terms of  $\gamma'$  and the last
   // column of  $\mathbf{U}'$ 
6:  $s \leftarrow 1$ 
7: if  $\gamma' = -\gamma$  then
8:    $s \leftarrow -1$ 
9: end if
10: for  $i = 1$  to  $d$  do
11:    $m_i \leftarrow s \cdot u'_{i, d+1}$ 
12: end for

```

---

**Remark 7.** Classically [29], [50], [51] Kannan’s idea is used in a much finer way, with several calls to lattice basis reduction, giving good theoretical guarantees of quality. We have however found that this rough (but efficient, as we only reduce one basis) version is sufficient for the problem at hand.

## V. LATTICE-BASED FILTER DESIGN

We start by discretizing Problem 2: we pick  $\ell (\geq n + 1)$  points  $\omega_0, \dots, \omega_{\ell-1}$  in  $\Omega$  and search for an approximation  $P(\omega) = \sum_{k=0}^n m_k \varphi_k(\omega)$ , where  $m_k \in I_b$  for  $k = 0, \dots, n$ , such that the vectors

$$\begin{pmatrix} W(\omega_0) \sum_{k=0}^n m_k \varphi_k(\omega_0) \\ W(\omega_1) \sum_{k=0}^n m_k \varphi_k(\omega_1) \\ \vdots \\ W(\omega_{\ell-1}) \sum_{k=0}^n m_k \varphi_k(\omega_{\ell-1}) \end{pmatrix} \text{ and } \begin{pmatrix} W(\omega_0) D(\omega_0) \\ W(\omega_1) D(\omega_1) \\ \vdots \\ W(\omega_{\ell-1}) D(\omega_{\ell-1}) \end{pmatrix}$$

are as close as possible with respect to  $\|\cdot\|_\infty$ . In other words, we need to find  $(m_0, \dots, m_n) \in I_b^{n+1}$  such that

$$m_0 \underbrace{\begin{pmatrix} W(\omega_0) \varphi_0(\omega_0) \\ W(\omega_1) \varphi_0(\omega_1) \\ \vdots \\ W(\omega_{\ell-1}) \varphi_0(\omega_{\ell-1}) \end{pmatrix}}_{\mathbf{b}_0} + \dots + m_n \underbrace{\begin{pmatrix} W(\omega_0) \varphi_n(\omega_0) \\ W(\omega_1) \varphi_n(\omega_1) \\ \vdots \\ W(\omega_{\ell-1}) \varphi_n(\omega_{\ell-1}) \end{pmatrix}}_{\mathbf{b}_n}$$

and  $\mathbf{t} = (W(\omega_0)D(\omega_0) \cdots W(\omega_{\ell-1})D(\omega_{\ell-1}))^t$  are as close as possible to each other, i.e.,

$$\text{minimize } \|m_0 \mathbf{b}_0 + \dots + m_n \mathbf{b}_n - \mathbf{t}\|_\infty, \quad (3)$$

which is a  $\text{CVP}_\infty$  instance, the Euclidean lattice under consideration being  $\mathbb{Z}\mathbf{b}_0 + \dots + \mathbb{Z}\mathbf{b}_n \subset \mathbb{R}^\ell$ . We approximately solve (3) in two steps, that we will present in more details in Subsection V-A:

- the existence of well-tuned, practical approximation algorithms in the Euclidean setting leads us to (approximately) solve the  $\|\cdot\|_2$  version of (3) instead:

$$\text{minimize } \|m_0 \mathbf{b}_0 + \dots + m_n \mathbf{b}_n - \mathbf{t}\|_2, \quad (4)$$

- using combinations of the short vectors computed during the first step, we “turn” around the approximate solution of the  $\text{CVP}_2$  instance (4) in order to improve the approximate solution of the  $\text{CVP}_\infty$  instance (3). In fact, as the reader will see in V-A, we use this vicinity search to directly improve our solution to Problem 2.

Note that working with (4) instead of directly solving (3) has the following important consequence. Since in  $\mathbb{R}^\ell$ ,  $\|\cdot\|_\infty \leq \|\cdot\|_2 \leq \sqrt{\ell} \|\cdot\|_\infty$ , we expect in general that small vectors with respect to the  $\|\cdot\|_2$  norm are also small when considering  $\|\cdot\|_\infty$ . Taking  $\ell$  close to the minimal value  $n$  hence helps [21].

**Remark 8.** It also proves very useful to approximate  $P^*$ , the minimax approximation, instead of  $D$ : in this case, we consider  $\mathbf{t} = (W(\omega_0)P^*(\omega_0) \cdots W(\omega_{\ell-1})P^*(\omega_{\ell-1}))^t$  in (3) and (4).

We have presented in Section III three families of discretization points:

- the alternating extrema of the minimax error function, cf. III-A, and approximate Fekete points, cf. III-C. We will use them to approximate both  $D$  and  $P^*$ ;
- the “zeros” of the minimax error function, cf. III-B. Some of these points may belong to  $[0, \pi] \setminus \Omega$ , hence we can use this family only to approximate  $P^*$ .

Eventually, we summarize our method in the form of a pseudo-code listing in Section V-B.

### A. Solving the CVP problem and a refinement trick

We saw in Section IV that (4) is an NP-hard problem, but argued that we can obtain an approximate solution quite efficiently by using a version of Kannan’s embedding technique. Applying Algorithm 3 to (4) determines:

- a lattice vector  $\sum_{k=0}^n m_k \mathbf{b}_k$  that is close to  $\mathbf{t}$  with respect to  $\|\cdot\|_2$  and hopefully close to  $\mathbf{t}$  with respect to  $\|\cdot\|_\infty$ ,

- but also a reduced basis  $(\mathbf{c}_0, \dots, \mathbf{c}_n)$  of  $(\mathbf{b}_0, \dots, \mathbf{b}_n)$  (by applying the LLL, BKZ or HKZ algorithm and retrieving the change of basis matrix  $\mathbf{U}$ ).

Since the  $\mathbf{c}_i$  vectors are usually short with respect to the  $\|\cdot\|_2$  norm (and consequently with the  $\|\cdot\|_\infty$  norm as well), we can use them to search in the vicinity of  $\sum_{k=0}^n m_k \mathbf{b}_k$  with the goal of potentially improving the quality of our result with respect to Problem 2.

In the function approximation setting, this idea is described in [57, p. 128]. Here, it translates to the following strategy:

---

**Algorithm 4** Vicinity search
 

---

**Input:** vector of discretized coefficients  $\mathbf{m} \in \mathbb{Z}^{n+1}$ , ideal response  $D$ , weight  $W$ , basis functions  $(\varphi_k)_{0 \leq k \leq n}$ , change of basis matrix  $\mathbf{U} \in \mathbb{Z}^{(n+1) \times (n+1)}$ .

**Output:** new vector of coefficients  $\mathbf{m}' \in \mathbb{Z}^{n+1}$ .

```

// Get initial approximation error & coefficients
1:  $E_{\min} \leftarrow \|W(\omega) (\sum_{k=0}^n m_k \varphi_k(\omega) - D(\omega))\|_{\Omega, \infty}$ 
2:  $\mathbf{m}' \leftarrow \mathbf{m}$ 
// Try to improve this initial approximation by selecting
// two directions  $d_0$  and  $d_1$  in which to search
3: for  $d_0 = 0$  to 8 do
4:   for  $d_1 = d_0 + 1$  to  $n$  do
5:      $\mathbf{u} \leftarrow \mathbf{U}_{0:n, d_0}, \mathbf{v} \leftarrow \mathbf{U}_{0:n, d_1}$ 
// Update coefficients in these directions if it leads
// to smaller approximation errors
6:     for  $\varepsilon_0, \varepsilon_1 \in \{0, \pm 1\}$  do
7:        $\mathbf{m}'' \leftarrow \mathbf{m}$ 
8:       for  $i = 0$  to  $n$  do
9:          $m''_i \leftarrow m''_i + u_i \varepsilon_0 + v_i \varepsilon_1$ 
10:      end for
11:       $E \leftarrow \|W(\omega) (\sum_{k=0}^n m''_k \varphi_k(\omega) - D(\omega))\|_{\Omega, \infty}$ 
12:      if  $E < E_{\min}$  then
13:         $E_{\min} \leftarrow E$ 
14:         $\mathbf{m}' \leftarrow \mathbf{m}''$ 
15:      end if
16:    end for
17:  end for
18: end for

```

---

We limit the value of  $d_0$  to a small constant (here to eight) in order to keep the execution time reasonable (line 11 takes  $O(n^2)$  operations, and lines 3, 4 & 6 tell us that we execute it  $O(n)$  times), but also because in practice we did not notice any significantly improved results by taking a larger search space. Increasing the number of search directions to three also helps, but we found the computational cost usually outweighs the improvements.

**Remark 9.** *The lattice basis reduction approaches we use provide results that are hopefully close to the actual CVP solution, but they are not necessarily optimal. And yet, there are two quite encouraging points:*

- the experiments in [40, §4.4.1] show that the output reduced basis is actually of excellent quality;
- the theoretical estimate in the same text [40, §4.4.3] supports our practical approach.

## B. A pseudo-code synthesis of our approach

We can synthesize our whole approach as follows:

---

**Algorithm 5** Lattice-based finite wordlength coefficient design
 

---

**Input:** degree  $n$ ,  $\ell$ -point discretization  $\{\omega_0, \dots, \omega_{\ell-1}\}$  of  $\Omega$  ( $\ell \geq n+1$ ), scaling factor  $s$ , minimax response  $P^*$ , weight  $W$ , a lattice basis reduction algorithm LBR  $\in \{\text{LLL, BKZ, HKZ}\}$ .

**Output:** fixed-point filter coefficients  $h[k], k = 0, \dots, 2n$ .

```

// Construct the appropriate basis functions  $\varphi_0, \dots, \varphi_n$ 
1:  $\varphi_0(\omega) \leftarrow \frac{1}{s}$ 
2: for  $k = 1$  to  $n$  do
3:    $\varphi_k(\omega) \leftarrow \frac{2 \cos(k\omega)}{s}$ 
4: end for
// Construct the lattice basis vectors
5: for  $i = 0$  to  $n$  do
6:    $\mathbf{b}_i \leftarrow (W(\omega_0) \varphi_i(\omega_0) \cdots W(\omega_{\ell-1}) \varphi_i(\omega_{\ell-1}))^t$ 
7: end for
// Construct our two possible target vectors
8:  $\mathbf{t} \leftarrow (W(\omega_0) P^*(\omega_0) \cdots W(\omega_{\ell-1}) P^*(\omega_{\ell-1}))^t$ 
9:  $\mathbf{t}' \leftarrow (W(\omega_0) D(\omega_0) \cdots W(\omega_{\ell-1}) D(\omega_{\ell-1}))^t$ 
// Find approximate CVP solutions using Algorithm 3
10:  $(\mathbf{m}, \mathbf{U}) \leftarrow \text{KannanEmbedding}(\mathbf{B}, \mathbf{t}, \text{LBR})$ 
11:  $(\mathbf{m}', \mathbf{U}') \leftarrow \text{KannanEmbedding}(\mathbf{B}, \mathbf{t}', \text{LBR})$ 
// Use the vicinity search of Algorithm 4 to improve the
// quality of the solution
12:  $\mathbf{m} \leftarrow \text{VicinitySearch}(\mathbf{m}, D, W, (\varphi_k)_{0 \leq k \leq n}, \mathbf{U})$ 
13:  $\mathbf{m}' \leftarrow \text{VicinitySearch}(\mathbf{m}', D, W, (\varphi_k)_{0 \leq k \leq n}, \mathbf{U}')$ 
14:  $E_1 \leftarrow \|W(\omega) (\sum_{k=0}^n m_k \varphi_k(\omega) - D(\omega))\|_{\Omega, \infty}$ 
15:  $E_2 \leftarrow \|W(\omega) (\sum_{k=0}^n m'_k \varphi_k(\omega) - D(\omega))\|_{\Omega, \infty}$ 
16: if  $E_2 < E_1$  then
17:    $\mathbf{m} \leftarrow \mathbf{m}'$ 
18: end if
// Retrieve the corresponding finite wordlength coefficients
19:  $h[n] \leftarrow \frac{m_0}{s}$ 
20: for  $k = 0$  to  $n-1$  do
21:    $h[k] \leftarrow h[2n-k] \leftarrow \frac{m_{n-k}}{s}$ 
22: end for

```

---

**Remark 10.** *Lines 10 and 11 will generate the same reduced basis for  $\mathbf{B}$ , so in practice we combine them to avoid any re-computations.*

**Remark 11.** *We call Algorithm 5 with the two discretization choices III-A and III-C. Regarding the discretization choice III-B, we call a reduced version of Algorithm 5: we don't execute lines 9, 11, 13 and 15 that all deal with the ideal function  $D$ . We keep the best result among the results returned from these calls.*

## VI. EXPERIMENTAL RESULTS

To illustrate the effectiveness of our method, we first compare it to the telescoping rounding approach introduced in [20], the error spectrum shaping method from [18] and the tree search approach of [4], based on least squares error optimization. To our knowledge, these<sup>6</sup> are the most efficient quasi-optimal

<sup>6</sup>Unfortunately, we were not able to make an accurate comparison with the method from [19].



TABLE II  
FILTER SPECIFICATIONS CONSIDERED IN [20].

Filter	Bands	$D(\omega)$	$W(\omega)$
A	$[0, 0.4\pi]$	1	1
	$[0.5\pi, \pi]$	0	1
B	$[0, 0.4\pi]$	1	1
	$[0.5\pi, \pi]$	0	10
C	$[0, 0.24\pi]$	1	1
	$[0.4\pi, 0.68\pi]$	0	1
	$[0.84\pi, \pi]$	1	1
D	$[0, 0.24\pi]$	1	1
	$[0.4\pi, 0.68\pi]$	0	10
	$[0.84\pi, \pi]$	1	1
E	$[0.02\pi, 0.42\pi]$	1	1
	$[0.52\pi, 0.98\pi]$	0	1

methods that explicitly treat the finite wordlength direct-form FIR design problem. We conclude this section with an account on the practical computation time of our code.

The following batch of tests were executed on an Intel i7-3687U CPU with a 64-bit Linux-based system and a g++ version 5.3.0 C++ compiler. As explained in Remark 10, all three discretization approaches presented in Section III are used, with the best result chosen in the end. On some of the examples, we will consider the passband peak to peak ripple  $20 \log_{10} \left( \frac{1+\delta_p}{1-\delta_p} \right)$  and passband and stopband attenuations  $-20 \log_{10} \delta_p$  and  $-20 \log_{10} \delta_s$  to measure the quality of our outputs, where  $\delta_p$  and  $\delta_s$  correspond to the unweighted passband and stopband approximation errors.

#### A. Kodek and Krisper's telescoping rounding [20]

We start with the type I FIR filter specifications given in [20, Table 1], which we reproduce for convenience, in Table II. When referring to a particular design instance, we give the specification letter, filter length and scaling factor. For example, A35/8 denotes a design problem adhering to specification A, of length  $N = 35$  (meaning a degree  $n = 17$  approximation),  $b = 8$  bits used to store the filter coefficients (sign bit included) and an  $s = 2^{b-1}$  scaling factor.

The results are highlighted in Table III:

- the first column lists the problem specification;
- the second one gives the minimax error  $E^*$  computed using the Parks-McClellan algorithm.

All remaining columns represent approximation errors of the transfer function for various coefficient quantization strategies:

- the third column presents optimal (or best known) finite wordlength errors  $E_{\text{opt}}$ : the errors for the length 125 filters are not proven to be optimal whereas the other ten errors are. Apart from the filters marked with the ‡, they are obtained using an MILP solver (with the values obtained in [20, Table 2]). In case of the problems marked with †, note that the exact solver was stopped after a certain time limit. We also indicate the methods attaining this optimal value, with M = (time-limited) MILP, L = lattice based (this paper), T = two coefficient telescoping rounding approach [20, Sec. 4].
- column four lists the errors obtained by simply rounding the real-valued minimax coefficients to their closest values in the imposed quantization format;

- the fifth column gives the best errors from using the two coefficient telescoping rounding approach of [20, Sec. 4];
- the last three columns show the lattice-based quantization errors when choosing the LLL, BKZ and HKZ basis reduction option, respectively, when applying Algorithm 5. A default block size of 8 was used when calling BKZ.

For the last four columns, values reported in bold are the best out of the telescoping rounding approach and our approach.

We notice that lattice-based quantization gives results which are optimal (or the best known ones) in eight cases out of fifteen and the other seven cases are very close to the optimal ones. It outperforms telescoping rounding in twelve cases out of fifteen and yields the same (optimal) result in a thirteenth case: the use of the LLL option gives a better result in twelve cases and an identical (optimal) result in a thirteenth case, the use of the BKZ option gives a better result in eleven cases and the use of the HKZ option in nine cases. We remark, in particular, the good behavior of the LLL algorithm in all 15 test cases, further emphasizing the idea that the lattice bases we use are close to being reduced. Our approach seems to work particularly well when the gap between the minimax error and the naive rounding error is significant, which is the case where one is most interested in improvements over the naive rounding filter. Eventually, note that, in the C125/21 and D125/22 cases, our approach returns (in less than 8 seconds, see Subsection VI-D) results that are better than the ones provided by (time-limited) MILP tools.

#### B. Nielsen's error spectrum shaping approach [18]

Consider now the finite wordlength low-pass filter specification [18, Sec. V] with passband  $[0, 0.4\pi]$ , stopband  $[0.6\pi, \pi]$  and maximum allowable attenuations of  $-90$  dB and  $-58.8$  dB in the stopband and passband, respectively. The corresponding weighting function is

$$W(\omega) = \begin{cases} 1, & \omega \in [0, 0.4\pi], \\ 10^{\frac{90-58.8}{20}}, & \omega \in [0.6\pi, \pi]. \end{cases}$$

The error spectrum shaping approach introduced in [18] and used by MATLAB<sup>TM</sup>'s fixed-point filter design tools gives a result with  $N = 69$ , 16-bit coefficients, stopband and passband attenuations of  $-97.7$  dB and  $-67.89$  dB, respectively [18, Table I]. Using our LLL-based routine, which took about 1.5 seconds to execute, we were able to obtain a smaller filter with  $N = 65$ , 15-bit coefficients and corresponding attenuations of  $-91.05$  dB and  $-61.27$  dB (see Table IV).

#### C. Lim and Parker's LMS criterion-based tree search approach [4]

Similarly, we can take the high-pass specification from Example 1 in [4] consisting of a  $[0, 0.74\pi]$  stopband,  $[0.8\pi, \pi]$  passband, 17 bit wordlength (including sign bit) coefficients, stopband attenuation  $\leq -80$  dB and passband peak to peak ripple  $\leq 0.1$  dB (passband attenuation  $-44.79$  dB). The result they obtain has degree  $n = 60$ , passband ripple 0.08 dB and stopband attenuation  $-80.3$  dB, leading to the weight function

$$W(\omega) = \begin{cases} 10^{\frac{80-44.79}{20}}, & \omega \in [0, 0.74\pi], \\ 1, & \omega \in [0.8\pi, \pi]. \end{cases}$$

TABLE III  
QUANTIZATION ERROR COMPARISON FOR THE FILTER SPECIFICATIONS GIVEN IN [20].

Filter	Minimax $E^*$	Best known finite wordlength $E_{\text{opt}}$	Naive rounding $E_{\text{naive}}$	Telescoping $E_{\text{tel}}$	LLL reduction $E_{\text{LLL}}$	BKZ reduction $E_{\text{BKZ}}$	HKZ reduction $E_{\text{HKZ}}$
A35/8	0.01595	0.02983 (M,L)	0.03266	0.03266	<b>0.02983</b>	<b>0.02983</b>	<b>0.02983</b>
A45/8	$7.132 \cdot 10^{-3}$	0.02962 (M,L)	0.03701	0.03186	<b>0.02962</b>	<b>0.02962</b>	<b>0.02962</b>
A125/21 <sup>†</sup>	$8.055 \cdot 10^{-6}$	$1.077 \cdot 10^{-5}$ (M)	$1.620 \cdot 10^{-5}$	$1.179 \cdot 10^{-5}$	<b><math>1.161 \cdot 10^{-5}</math></b>	$1.168 \cdot 10^{-5}$	$1.251 \cdot 10^{-5}$
B35/9	0.05275	0.07709 (M)	0.15879	<b>0.07854</b>	0.08205	0.08205	0.08205
B45/9	0.02111	0.05679 (M)	0.11719	0.06641	<b>0.06041</b>	<b>0.06041</b>	<b>0.06041</b>
B125/21 <sup>†</sup>	$2.499 \cdot 10^{-5}$	$2.959 \cdot 10^{-5}$ (M)	$6.198 \cdot 10^{-5}$	$3.293 \cdot 10^{-5}$	<b><math>3.243 \cdot 10^{-5}</math></b>	$3.344 \cdot 10^{-5}$	$3.344 \cdot 10^{-5}$
C35/8	$2.631 \cdot 10^{-3}$	0.01787 (M,T,L)	0.04687	<b>0.01787</b>	<b>0.01787</b>	0.01917	0.01917
C45/8	$6.709 \cdot 10^{-4}$	0.01609 (M,L)	0.03046	0.02103	<b>0.01609</b>	<b>0.01609</b>	0.02291
C125/21 <sup>‡</sup>	$1.278 \cdot 10^{-8}$	$1.564 \cdot 10^{-6}$ (L)	$8.203 \cdot 10^{-6}$	$2.1 \cdot 10^{-6}$	$1.606 \cdot 10^{-6}$	<b><math>1.564 \cdot 10^{-6}</math></b>	<b><math>1.564 \cdot 10^{-6}</math></b>
D35/9	0.01044	0.03252 (M,L)	0.12189	0.03368	<b>0.03252</b>	0.03291	0.03291
D45/9	$2.239 \cdot 10^{-3}$	0.02612 (M)	0.10898	0.02859	<b>0.02706</b>	0.02805	<b>0.02706</b>
D125/22 <sup>‡</sup>	$4.142 \cdot 10^{-8}$	$1.781 \cdot 10^{-6}$ (L)	$3.425 \cdot 10^{-5}$	$2.16 \cdot 10^{-6}$	$1.864 \cdot 10^{-6}$	<b><math>1.781 \cdot 10^{-6}</math></b>	$1.826 \cdot 10^{-6}$
E35/8	0.01761	0.03299 (M)	0.04692	0.03404	<b>0.03349</b>	<b>0.03349</b>	<b>0.03349</b>
E45/8	$6.543 \cdot 10^{-3}$	0.02887 (M,L)	0.03571	0.03403	0.03167	0.03094	<b>0.02887</b>
E125/21 <sup>†</sup>	$7.889 \cdot 10^{-6}$	$1.034 \cdot 10^{-5}$ (M)	$1.479 \cdot 10^{-5}$	<b><math>1.127 \cdot 10^{-5}</math></b>	$1.215 \cdot 10^{-5}$	$1.245 \cdot 10^{-5}$	$1.162 \cdot 10^{-5}$

TABLE IV

IMPULSE RESPONSE FOR THE DISCRETE COEFFICIENT FILTER MEETING THE SPECIFICATIONS OF [18, SEC. 5]  
 $h[n] = h[64 - n]$  FOR  $33 \leq n \leq 64$ .

Lowpass filter, $N = 65$			
Peak passband attenuation = $-61.27$ dB			
Peak stopband attenuation = $-91.05$ dB			
Impulse response $\times 2^{15}$			
$h[32] = 8470$	$h[21] = -365$	$h[10] = -11$	
$h[31] = 5206$	$h[20] = 129$	$h[9] = -41$	
$h[30] = -272$	$h[19] = 273$	$h[8] = 5$	
$h[29] = -1710$	$h[18] = -94$	$h[7] = 25$	
$h[28] = 256$	$h[17] = -200$	$h[6] = 0$	
$h[27] = 995$	$h[16] = 66$	$h[5] = -13$	
$h[26] = -231$	$h[15] = 145$	$h[4] = -2$	
$h[25] = -677$	$h[14] = -41$	$h[3] = 5$	
$h[24] = 199$	$h[13] = -100$	$h[2] = 1$	
$h[23] = 490$	$h[12] = 24$	$h[1] = -2$	
$h[22] = -165$	$h[11] = 67$	$h[0] = -1$	

By using LLL as the basis reduction algorithm, we obtain a similar filter of degree  $n = 60$ , with passband ripple 0.091 dB and stopband attenuation  $-80.03$  dB in under 9 seconds; the resulting coefficients are given in Table V.

#### D. Computation time in practice

Analyzing the runtime of Algorithm 5 is dependent on the basis reduction approach used for the Kannan embedding portion of the computation (lines 10 and 11) (see last columns of Table I) and the short vector-based vicinity search (lines 12 and 13), which takes  $O(n^3)$ . Actually, the lattice bases that appear in our problems are usually quite close to being reduced, making the Kannan embedding portion of the code much faster than the vicinity search (at least when using LLL and BKZ reductions). This can be seen for instance in Table VI, which breaks down the execution time of our code, applied to the three discretizations mentioned in Section III, on the examples of Table III:

- the first column lists the problem specification;
- the second one is the time required for finding the three discretizations;
- the following three columns show the computation time of the Euclidean lattice part of the code (lattice basis

TABLE V

IMPULSE RESPONSE FOR THE DISCRETE COEFFICIENT FILTER MEETING THE SPECIFICATIONS OF [3, EXAMPLE 1]  
 $h[n] = h[120 - n]$  FOR  $61 \leq n \leq 120$ .

Highpass filter, $N = 121$			
Passband peak to peak ripple = 0.091 dB			
Peak stopband attenuation = $-80.03$ dB			
Impulse response $\times 2^{16}$			
$h[60] = 14705$	$h[39] = -593$	$h[18] = -139$	
$h[59] = -13508$	$h[38] = 148$	$h[17] = 131$	
$h[58] = 10268$	$h[37] = 301$	$h[16] = -71$	
$h[57] = -5914$	$h[36] = -555$	$h[15] = -6$	
$h[56] = 1631$	$h[35] = 529$	$h[14] = 64$	
$h[55] = 1536$	$h[34] = -274$	$h[13] = -81$	
$h[54] = -3010$	$h[33] = -65$	$h[12] = 57$	
$h[53] = 2817$	$h[32] = 325$	$h[11] = -9$	
$h[52] = -1502$	$h[31] = -402$	$h[10] = -39$	
$h[51] = -135$	$h[30] = 287$	$h[9] = 67$	
$h[50] = 1358$	$h[29] = -58$	$h[8] = -68$	
$h[49] = -1749$	$h[28] = -167$	$h[7] = 46$	
$h[48] = 1308$	$h[27] = 289$	$h[6] = -14$	
$h[47] = -377$	$h[26] = -269$	$h[5] = -15$	
$h[46] = -559$	$h[25] = 136$	$h[4] = 32$	
$h[45] = 1095$	$h[24] = 35$	$h[3] = -35$	
$h[44] = -1063$	$h[23] = -163$	$h[2] = 28$	
$h[43] = 566$	$h[22] = 199$	$h[1] = -17$	
$h[42] = 107$	$h[21] = -141$	$h[0] = 8$	
$h[41] = -634$	$h[20] = 29$		
$h[40] = 804$	$h[19] = 80$		

reduction and approximate CVP solving, which correspond to lines 10 and 11 in Algorithm 5) when choosing the LLL, BKZ and HKZ basis reduction option, respectively;

- column six presents the vicinity search computation time, which correspond to lines 12 and 13 in Algorithm 5;
- the last column gives the total execution time of our code when choosing the LLL option. The remaining steps are executed fast enough so that the values in this column are just slightly larger than the sum of the values presented in columns 2, 3 and 6).

We used the `fp111` C++ library implementations [58] of the LLL, BKZ and HKZ algorithms.

One can note that our approach, when choosing the LLL option, is fast: for instance, it takes at most 8 seconds to compute a very good frequency response for the filters of length 125. Interestingly enough, the use of the BKZ option,

TABLE VI

TIMINGS (IN SECONDS) WHEN USING THE THREE DISCRETIZATIONS OF SECTION III. THE LAST COLUMN GIVES THE TOTAL RUNTIME WHEN LLL IS CHOSEN AS THE LATTICE BASIS REDUCTION ALGORITHM.

Filter	Point generation	LLL	BKZ	HKZ	Vicinity search	Total (LLL)
A35/8	0.048	0.026	0.028	0.028	0.231	0.325
A45/8	0.047	0.033	0.035	0.035	0.474	0.571
A125/21	0.112	0.291	0.348	0.471	7.686	8.185
B35/9	0.112	0.018	0.019	0.027	0.243	0.401
B45/9	0.102	0.033	0.035	0.045	0.468	0.676
B125/22	0.191	0.296	0.342	1.194	7.746	8.314
C35/8	0.094	0.022	0.021	0.021	0.225	0.396
C45/8	0.107	0.032	0.036	0.037	0.414	0.611
C125/21	0.319	0.455	0.671	303.98	6.366	7.278
D35/9	0.162	0.019	0.021	0.021	0.211	0.411
D45/9	0.096	0.031	0.037	0.037	0.411	0.595
D125/22	0.251	0.435	0.784	335.38	6.240	7.042
E35/8	0.063	0.017	0.018	0.021	0.211	0.301
E45/8	0.089	0.018	0.031	0.031	0.420	0.574
E125/21	0.236	0.261	0.312	1.935	1.184	7.712

TABLE VII

HIGH DEGREE TYPE I FIR SPECIFICATIONS

Filter	Bands	$D(\omega)$	$W(\omega)$
F	$[0, 0.2\pi]$	0	10
	$[0.205\pi, \pi]$	1	1
G	$[0, 0.4\pi]$	1	1
	$[0.405\pi, 0.7\pi]$	0	10
	$[0.705\pi, \pi]$	1	1
H	$[0, 0.45\pi]$	0	10
	$[0.47\pi, \pi]$	1	1

which computes better reduced bases, leads to a very small computing time penalty. The use of the BKZ option allows us to compute the best, as of today, frequency responses for C125 and D125 in less than 8 seconds. Another remarkable fact is the excellent behavior of our approach when choosing the HKZ option: it is very fast in the cases of A125, B125 and E125 and offers a reasonable execution time in the cases of C125, D125 whereas the corresponding lattices have a dimension equal to, at least, 63, a size which usually means a very difficult challenge for an HKZ reduction attempt.

**Remark 12.** *To further illustrate its robustness, let's mention that our Euclidean lattice-based routine scales well to much larger degrees, as illustrated in Table VIII, which adhere to the specifications from Table VII. The scaling factor used is again  $s = 2^{b-1}$ . Note that:*

- since the exhaustive enumeration performed by the vicinity search would be quite costly on such examples, we use a simple stochastic approach: take uniform random values of  $d_0, d_1, \varepsilon_0, \varepsilon_1$  inside Algorithm 4 for one minute, keeping the best result at the end;

TABLE VIII

HIGH DEGREE QUANTIZATION RESULTS

Filter	Minimax $E^*$	Naive rounding $E_{\text{naive}}$	LLL reduction $E_{\text{LLL}}$	Runtime in sec.
F1601/18	$8.64 \cdot 10^{-4}$	$3.41 \cdot 10^{-3}$	$1.20 \cdot 10^{-3}$	846.85
G1201/16	$4.78 \cdot 10^{-3}$	$1.21 \cdot 10^{-2}$	$5.79 \cdot 10^{-3}$	567.01
H501/16	$1.67 \cdot 10^{-4}$	$5.79 \cdot 10^{-3}$	$1.44 \cdot 10^{-3}$	379.11

- if the filter degree is too large, at some point the leading coefficients of the optimal quantization will become zero [59]. We have avoided this scenario by considering only very narrow transition bands, which correspond to a slow decrease in the size of the minimax filter coefficients.

## VII. CONCLUSION

We have developed a novel approach for designing machine-number coefficient FIR filters based on an idea previously introduced in [21], which transforms the design problem using the language of Euclidean lattices. The values obtained show that the method is extremely robust and competitive in practice. It frequently produces results which are close to optimal and/or beats other heuristic approaches.

There are several directions of research we are currently pursuing. The first is integrating this method into a FIR filter design suite for FPGA targets. The text [60] is a first attempt in this direction. IIR filter quantization using Euclidean lattices is another idea. There are several difficulties which have to be overcome in the rational IIR setting:

- nonlinearity of the transfer function;
- the approximation domain switches from  $\Omega$  to a subset of the unit circle;
- ensuring stability of the transfer function (control the position of poles).

## REFERENCES

- [1] Y. C. Lim and A. Constantinides, "New integer programming scheme for nonrecursive digital filter design," *Electronics Letters*, vol. 15, no. 25, pp. 812–813, Dec. 1979.
- [2] Y. C. Lim, S. Parker, and A. Constantinides, "Finite word length FIR filter design using integer programming over a discrete coefficient space," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 30, no. 4, pp. 661–664, Aug. 1982.
- [3] Y. C. Lim and S. Parker, "FIR filter design over a discrete powers-of-two coefficient space," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 31, no. 3, pp. 583–591, Jun. 1983.
- [4] —, "Discrete Coefficient Fir Digital Filter Design Based Upon an LMS Criteria," *IEEE Transactions on Circuits and Systems*, vol. 30, no. 10, pp. 723–739, Oct. 1983.
- [5] Y. C. Lim, R. Yang, D. Li, and J. Song, "Signed power-of-two term allocation scheme for the design of digital filters," *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*, vol. 46, no. 5, pp. 577–584, May 1999.
- [6] J. Yli-Kaakinen and T. Saramäki, "A Systematic Algorithm for the Design of Lattice Wave Digital Filters With Short-Coefficient Wordlength," *IEEE Transactions on Circuits and Systems*, vol. 54, no. 8, pp. 1838–1851, Aug. 2007.
- [7] J. Skaf and S. P. Boyd, "Filter Design With Low Complexity Coefficients," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3162–3169, Jul. 2008.
- [8] Y. Cao, K. Wang, W. Pei, Y. Liu, and Y. Zhang, "Design of high-order extrapolated impulse response FIR filters with signed powers-of-two coefficients," *Circuits, Systems, and Signal Processing*, vol. 30, no. 5, pp. 963–985, 2011.
- [9] E. A. da Silva, L. Lovisolio, A. J. Dutra, and P. S. Diniz, "FIR filter design based on successive approximation of vectors," *IEEE Transactions on Signal Processing*, vol. 62, no. 15, pp. 3833–3848, 2014.
- [10] J. Knowles and E. Olcayto, "Coefficient accuracy and digital filter response," *IEEE Transactions on Circuit Theory*, vol. 15, no. 1, pp. 31–41, 1968.
- [11] D. M. Kodek, "Design of optimal finite wordlength FIR digital filters using integer programming techniques," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 28, no. 3, pp. 304–308, Jun. 1980.
- [12] T. W. Parks and J. H. McClellan, "Chebyshev Approximation for Nonrecursive Digital Filters with Linear Phase," *IEEE Transactions on Circuit Theory*, vol. 19, no. 2, pp. 189–194, 1972.



- [13] D. M. Kodek, "Performance limit of finite wordlength FIR digital filters," *IEEE Transactions on Signal Processing*, vol. 53, no. 7, pp. 2462–2469, Jul. 2005.
- [14] D. M. Kodek and K. Steiglitz, "Comparison of optimal and local search methods for designing finite wordlength fir digital filters," *IEEE Transactions on Circuits and Systems*, vol. 28, no. 1, pp. 28–32, 1981.
- [15] O. Gustafsson and L. Wanhammar, "Design of linear-phase FIR filters combining subexpression sharing with MILP," in *The 2002 45th Midwest Symposium on Circuits and Systems, 2002. MWSCAS-2002.*, vol. 3, Aug. 2002, pp. 9–12.
- [16] D. Shi and Y. J. Yu, "Design of Linear Phase FIR Filters With High Probability of Achieving Minimum Number of Adders," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 58, no. 1, pp. 126–136, Jan. 2011.
- [17] D. M. Kodek, "LLL Algorithm and the Optimal Finite Wordlength FIR Design," *IEEE Transactions on Signal Processing*, vol. 60, no. 3, pp. 1493–1498, Mar. 2012.
- [18] J. Nielsen, "Design of linear-phase direct-form FIR digital filters with quantized coefficients using error spectrum shaping," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 37, no. 7, pp. 1020–1026, Jul. 1989.
- [19] G. Evangelista, "Design of optimum high-order finite-wordlength digital FIR filters with linear phase," *Signal Processing*, vol. 82, no. 2, pp. 187–194, 2002.
- [20] D. M. Kodek and M. Krisper, "Telescoping rounding for suboptimal finite wordlength FIR digital filter design," *Digital Signal Processing*, vol. 15, no. 6, pp. 522–535, 2005.
- [21] N. Brisebarre and S. Chevillard, "Efficient polynomial  $L^\infty$ -approximations," in *18th IEEE Symposium on Computer Arithmetic, 2007. ARITH '07*, Jun. 2007, pp. 169–176.
- [22] L. P. Bos and N. Levenberg, "On the calculation of approximate Fekete points: the univariate case," *Electronic Transactions on Numerical Analysis*, vol. 30, pp. 377–397, 2008.
- [23] J.-P. Calvi and N. Levenberg, "Uniform approximation by discrete least squares polynomials," *Journal of Approximation Theory*, vol. 152, no. 1, pp. 82–100, May 2008.
- [24] A. Sommariva and M. Vianello, "Computing approximate Fekete points by QR factorizations of Vandermonde matrices," *Computers & Mathematics with Applications*, vol. 57, no. 8, pp. 1324–1336, Apr. 2009.
- [25] —, "Approximate Fekete points for weighted polynomial interpolation," *Electronic Transactions on Numerical Analysis*, vol. 37, pp. 1–22, 2010.
- [26] L. P. Bos, J.-P. Calvi, N. Levenberg, A. Sommariva, and M. Vianello, "Geometric weakly admissible meshes, discrete least squares approximations and approximate Fekete points," *Mathematics of Computation*, vol. 80, no. 275, pp. 1623–1638, 2011.
- [27] S. De Marchi, F. Piazzon, A. Sommariva, and M. Vianello, "Polynomial Meshes: Computation and Approximation," in *Proceedings of the 15th International Conference on Computational and Mathematical Methods in Science and Engineering*, 2015, pp. 414–425.
- [28] P. Q. Nguyen and B. Vallée, Eds., *The LLL Algorithm - Survey and Applications*, ser. Information Security and Cryptography. Springer, 2010. [Online]. Available: <http://dx.doi.org/10.1007/978-3-642-02295-1>
- [29] R. Kannan, "Minkowski's convex body theorem and integer programming," *Mathematics of Operations Research*, vol. 12, no. 3, pp. 415–440, Aug. 1987.
- [30] A. V. Oppenheim and R. W. Schaffer, *Discrete-Time Signal Processing*, ser. Prentice-Hall Signal Processing Series. Prentice Hall, 2010.
- [31] A. Antoniou, *Digital Signal Processing: Signals, Systems, and Filters*. McGraw-Hill Education, 2005.
- [32] P. Prandoni and M. Vetterli, *Signal Processing for Communications*. Taylor & Francis, 2008. [Online]. Available: <http://www.sp4comm.org/>
- [33] E. W. Cheney, *Introduction to Approximation Theory*, ser. AMS Chelsea Publishing Series. AMS Chelsea Pub., 1982.
- [34] D. M. Kodek, "Design of optimal finite wordlength FIR digital filters," in *Proc. of the European Conference on Circuit Theory and Design, ECCTD '99.*, vol. 1, Aug. 1999, pp. 401–404.
- [35] Y. Lim, "Design of discrete-coefficient-value linear phase FIR filters with optimum normalized peak ripple magnitude," *IEEE Transactions on Circuits and Systems*, vol. 37, no. 12, pp. 1480–1486, Dec. 1990.
- [36] L. N. Trefethen, *Approximation Theory and Approximation Practice*. Society for Industrial and Applied Mathematics, 2013.
- [37] S.-I. Filip, "A robust and scalable implementation of the Parks-McClellan algorithm for designing FIR filters," *ACM Transactions on Mathematical Software*, vol. 43, no. 1, Aug. 2016. [Online]. Available: <https://hal.inria.fr/hal-01136005>
- [38] A. Çivril and M. Magdon-Ismael, "On selecting a maximum volume sub-matrix of a matrix and related problems," *Theoretical Computer Science*, vol. 410, no. 47–49, pp. 4801–4811, Nov. 2009.
- [39] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen, *LAPACK Users' guide*. SIAM, 1999, vol. 9.
- [40] S.-I. Filip, "Robust tools for weighted Chebyshev approximation and applications to digital filter design," Ph.D. dissertation, École Normale Supérieure de Lyon, 2016. [Online]. Available: <https://tel.archives-ouvertes.fr/tel-01447081/>
- [41] M. Embree and L. N. Trefethen, "Green's Functions for Multiply Connected Domains via Conformal Mapping," *SIAM Review*, vol. 41, no. 4, pp. 745–761, Dec. 1999.
- [42] J. Shen and G. Strang, "The asymptotics of optimal (equiripple) filters," *IEEE Trans. on Signal Processing*, vol. 47, pp. 1087–1098, 1999.
- [43] J. Shen, G. Strang, and A. J. Wathen, "The potential theory of several intervals and its applications," *Appl. Math. Opt.*, vol. 44, pp. 67–85, 2001.
- [44] P. M. Gruber and C. G. Lekkerkerker, *Geometry of numbers*, 2nd ed., ser. North-Holland Mathematical Library. Amsterdam: North-Holland Publishing Co., 1987, vol. 37.
- [45] L. Lovász, *An algorithmic theory of numbers, graphs and convexity*, ser. CBMS-NSF Regional Conference Series in Applied Mathematics. Philadelphia, PA: Society for Industrial and Applied Mathematics (SIAM), 1986, vol. 50.
- [46] J. H. Conway and N. J. A. Sloane, *Sphere Packings, Lattices and Groups*. New York: Springer-Verlag, 1988.
- [47] D. Micciancio and S. Goldwasser, *Complexity of Lattice Problems: a cryptographic perspective*, ser. The Kluwer International Series in Engineering and Computer Science. Boston, Massachusetts: Kluwer Academic Publishers, Mar. 2002, vol. 671.
- [48] T. H. I. U. of Crystallography, *Space-group symmetry*, ser. International Tables for Crystallography. Springer Netherlands, 2002, vol. A.
- [49] P. van Emde Boas, "Another NP-complete problem and the complexity of computing short vectors in a lattice," University of Amsterdam, Department of Mathematics, Netherlands, Tech. Rep. 8104, 1981. [Online]. Available: <https://staff.fnwi.uva.nl/p.vanemdeboas/>
- [50] R. Kannan, "Algorithmic geometry of numbers," *Annual Review of Computer Science*, vol. 2, pp. 231–267, 1987.
- [51] C. K. Dubey and T. Holenstein, "Approximating the closest vector problem using an approximate shortest vector oracle," in *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques - 14th International Workshop, APPROX 2011, and 15th International Workshop, RANDOM 2011, Princeton, NJ, USA, August 17-19, 2011. Proceedings*, 2011, pp. 184–193.
- [52] A. K. Lenstra, H. W. Lenstra, Jr., and L. Lovász, "Factoring polynomials with rational coefficients," *Math. Ann.*, vol. 261, pp. 515–534, 1982.
- [53] C. P. Schnorr, "A hierarchy of polynomial lattice basis reduction algorithms," *Theoretical Computer Science*, vol. 53, pp. 201–224, 1987.
- [54] G. Hanrot, X. Pujol, and D. Stehlé, "Analyzing blockwise lattice algorithms using dynamical systems," in *Advances in Cryptology - CRYPTO 2011 - 31st Annual Cryptology Conference, Santa Barbara, CA, USA, August 14-18, 2011. Proceedings*, 2011, pp. 447–464.
- [55] R. Kannan, "Improved algorithms for integer programming and related lattice problems," in *Proceedings of the 15th Symposium on the Theory of Computing (STOC 1983)*. ACM, 1983, pp. 99–108.
- [56] D. Aggarwal, D. Dadush, and N. Stephens-Davidowitz, "Solving the Closest Vector Problem in  $2^n$  Time – The Discrete Gaussian Strikes Again!" in *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*. IEEE, 2015, pp. 563–582.
- [57] S. Chevillard, "Évaluation efficace de fonctions numériques. outils et exemples," Ph.D. dissertation, École Normale Supérieure de Lyon, 2009. [Online]. Available: <http://www-sop.inria.fr/members/Sylvain.Chevillard/>
- [58] The FPLLL development team, "fpLLL, a lattice reduction library," 2016, available at <https://github.com/fplll/fplll>. [Online]. Available: <https://github.com/fplll/fplll>
- [59] D. M. Kodek, "Length limit of optimal finite wordlength FIR filters," *Digital Signal Processing*, vol. 25, no. 5, pp. 1798–1805, Sep. 2013.
- [60] N. Brisebarre, F. de Dinechin, S.-I. Filip, and M. Istoan, "Automatic generation of hardware FIR filters from a frequency domain specification," 2016. [Online]. Available: <https://hal.inria.fr/hal-01308377>