

# Architecture synthesis for linear time-invariant filters

Antoine Martinet

Master 2 Internship report

at CITI lab,

under the supervision of

Florent de Dinechin

2 February - 31 July,  
2015

## Contents

<b>Introduction</b>	<b>3</b>
<b>1 Tools for expressing filters and signal processing</b>	<b>3</b>
1.1 Fundamentals about signal processing	3
1.1.1 Definition of a signal	3
1.1.2 $l^\infty$ norm	3
1.1.3 Linear Time Invariant Filters (LTI filters)	3
1.1.4 Impulse response	4
1.1.5 Worst-Case Peak Gain (WCPG) of a Filter: the maximum possible amplification	4
1.2 FIR and IIR: two filters families	4
1.3 Different realizations: how to compute the output of a filter?	5
1.3.1 Direct and transposed forms	5
1.3.2 State-space representation: a recurrence to define an infinite response	5
1.4 The Specialized Implicit Form (SIF): a unified representation	5
1.4.1 Definition	5
1.4.2 Workflow for filter implementation	7
<b>2 Architecture generation</b>	<b>8</b>
2.1 Implementation: what can be done in practise?	8
2.1.1 The FloPoCo Framework: computing just right	8
2.1.2 KCM multipliers: an architecture to efficiently compute products on FPGAs	8
2.1.3 SOPCs: exploiting full parallel architectures	8
2.1.4 Precision concerns	8
2.1.5 Size computation in finite precision	9
2.1.6 Error analysis	10
2.1.7 Worst-case peak gain computation	12
2.2 Architecture generation algorithm	12
2.3 Particular Forms, degenerated cases	12
2.3.1 ABCD Form: still achievable	12
2.3.2 When $n_x = 0$	14
2.4 Optimizations to do on this technology	14
2.4.1 Sparse matrices	14
2.4.2 Power of two coefficients	14
2.5 Implementation details	14
2.6 Filter specification interface	14
<b>3 Further work</b>	<b>15</b>
3.1 Static analysis on coefficients	15
3.2 Sub-filter detection	15
3.3 Precision calculations improvement	15
3.4 File format re-specification	15
<b>Conclusion</b>	<b>16</b>

---

## Introduction

Filters are nowadays essential tools for designing responsive systems. In signal processing, filters are used in radio signal coding and decoding, image and sound processing, and so on. In addition, they are also used in many control applications.

For most applications, filters can be computed in software. However, for performance reasons hardware implementations are needed in many applications. There is an interest in FPGAs implementations because they can help for prototyping a software defined radio (SDR).

SDR is the main research subject of the SOCRATE team, at the CITI lab, where this work has been conducted. The goal of this internship was to design a parametric architecture generator for Linear Time Invariant (LTI) filters. This work has been integrated to FloPoCo tool, whose purpose is to generate architecture cores computing just right.

This report first presents the context of LTI filters, SIF and their implementation in part 1. Part 2 will then describe the details of the implementation brought in arbitrary precision. Finally, an overview of the future work will be presented in part 3.

## 1 Tools for expressing filters and signal processing

### 1.1 Fundamentals about signal processing

#### 1.1.1 Definition of a signal

Lopez's PhD [6] gives a good presentation of the state of the art. Many notations and definitions are kept from this PhD.

Note that a bold symbol (for example  $\mathbf{h}$ ) denotes a matrix, that may be a vector. On the contrary, a normal symbol (for example  $\varepsilon_{t_1}(k)$ ) denotes a single value

**Definition 1.** (*Signal*) Generally, a signal is a temporal variable, which takes a value from  $\mathbb{R}$  at each time  $t$ .  $x(t)$  denotes the value of the signal  $x$  at the instant  $t$ . When dealing with discrete time events, the time will be represented by  $k$ . The notation will then be  $x(k)$ , which is said to be a sample.  $\{x(k)\}_{k \geq 0}$  denotes all the values possible for the signal  $x$ . The rest of this report addresses vectors of signals  $\mathbf{x}$ , where  $\mathbf{x}(k) \in \mathbb{R}^n$

#### 1.1.2 $\ell^\infty$ norm

**Definition 2.** ( $\ell^\infty$  norm) The  $\ell^\infty$  norm of a scalar signal  $x$ , denoted  $\|x\|_{\ell^\infty}$ , is the smallest upper bound among all values (absolute values) possible for the signal  $x$ , that is:

$$\|x\|_{\ell^\infty} = \sup_{k \in \mathbb{N}} |x(k)| \quad (1.1.1)$$

#### 1.1.3 Linear Time Invariant Filters (LTI filters)

A filter, denoted by its transfer function  $\mathcal{H}$ , is an application which transforms a signal vector  $\mathbf{u}$  (with  $\dim(\mathbf{u}) = n_u$ ) into a signal vector  $\mathbf{y} = \mathcal{H}(\mathbf{u})$ , of size  $\dim(\mathbf{y}) = n_y$ . The case where  $n_u = n_y = 1$  is referred as Single Input Single Output (SISO) filters. Other cases are referred as Multiple Input Multiple Output (MIMO) filters.

**Definition 3.** (*Linear Time Invariant Filter*)

*Linearity:*

$$\mathcal{H}(\alpha \cdot \mathbf{u}_1 + \beta \cdot \mathbf{u}_2) = \alpha \cdot \mathcal{H}(\mathbf{u}_1) + \beta \cdot \mathcal{H}(\mathbf{u}_2)$$

*Time invariance:*

$$\{\mathcal{H}(\mathbf{u})(k - k_0)\}_{k \geq 0} = \mathcal{H}(\{\mathbf{u}(k - k_0)\}_{k \geq 0})$$

### 1.1.4 Impulse response

**Definition 4.** (*Impulse Response*) A SISO filter may be defined by its impulse response, denoted  $h$ .  $h$  is the impulse response of  $H$  to the impulsion of Dirac. Indeed each input signal can be described as a sum of Dirac impulsions:

$$u = \sum_{l \geq 0} u(l) \delta_l$$

where  $\delta_l$  is a Dirac impulsion centered in  $l$ , that is:

$$\delta_l(k) = \begin{cases} 1 & \text{when } k = l \\ 0 & \text{else} \end{cases} \quad (1.1.2)$$

The linearity condition of  $\mathcal{H}$  implies:  $\mathcal{H}(u) = \sum_{l \geq 0} u(l) \mathcal{H}(\delta_l)$ . Time invariance gives:  $\mathcal{H}(\delta_l)(k) = h(k - l)$ . Then the computation from inputs to outputs takes this form:

$$y(k) = \sum_{l \geq 0} u(l) h(k - l) = \sum_{l=0}^k u(l) h(k - l)$$

This corresponds with the convolution product definition of  $u$  by  $h$ , denoted  $y = h * u$ . Dealing with MIMO filters,  $\mathbf{h}(k) \in \mathbb{R}^{n_y \times n_u}$  is the impulse response of  $\mathcal{H}$ .  $\mathbf{h}_{i,j}(k)$  is the response on the  $i$ th output to the Dirac impulsion on the  $j$ -th input. The precedent equation becomes:

$$y_i(k) = \sum_{j=1}^{n_u} \sum_{l \geq 0} u_j(l) h_{i,j}(k - l), \quad \forall 1 \leq i \leq n_y$$

### 1.1.5 Worst-Case Peak Gain (WCPG) of a Filter: the maximum possible amplification

**Definition 5.** (*Worst-Case Peak Gain*) The worst case peak gain is defined as the maximum amplification possible over all potential inputs through the filter.

$$\|\mathcal{H}\|_{wcpg} = \sup_{u \neq 0} \frac{\|h * u\|_{l^\infty}}{\|u\|_{l^\infty}}$$

with  $h$  the impulse response of  $\mathcal{H}$ ,  $u$  the input signal, and  $h * u$  the convolution product of  $h$  by  $u$  (output of the filter).

## 1.2 FIR and IIR: two filters families

There are two types of LTI filters: *Finite impulse response* (FIR) and *Infinite impulse response* (IIR) filters. Formally, the impulse response is finite when:

$$\exists n \in \mathbb{N} \forall k \geq n, h(k) = 0 \quad (1.2.1)$$

The smallest  $n$  verifying 1.2.1 is referred as the order of the filter. So a  $n$ -order FIR can be described by the following equation:

$$y(k) = \sum_{i=0}^n b_i u(k - i) \quad (1.2.2)$$

An IIR will be described as following:

$$y(k) = \sum_{i=0}^n b_i u(k - i) - \sum_{i=0}^n a_i y(k - i) \quad (1.2.3)$$

Here one can observe that the output at time  $k$  depends also on all previous  $n$  outputs (loopback). Also remark that a FIR can be seen as an IIR with  $\forall i \in [0, n], a_i = 0$ . The impulse response can then be deduced from 1.2.3 by resolving the recurrence relation:

$$h(k) = \begin{cases} 0 & \text{when } k < l \\ b_k - \sum_{l=1}^n a_l h(k-l) & \text{when } 0 \leq k \leq n \\ \sum_{l=1}^n a_l h(k-l) & \text{when } n < k \end{cases} \quad (1.2.4)$$

### 1.3 Different realizations: how to compute the output of a filter?

**Definition 6.** (realization) A realization can be defined as an algorithm describing how to compute outputs from inputs. However, a realization does not describes the details of basic operations (format, size, rounding, etc...)

It is important to know that all realizations of a filter are mathematically equivalent to each other (infinite precision). But in finite precision, rounding aspects have a huge impact on the correctness of results. Most of the time in this report, realizations will be referred as forms.

#### 1.3.1 Direct and transposed forms

Direct and transposed forms are classic realizations. A good description of these forms can be found in Lopez' and Hilaire's PhDs [6] [5]. The direct form has been implemented in the FoPoCo project, The hardware implementation of anIIR can be seen on figure 1.3.1

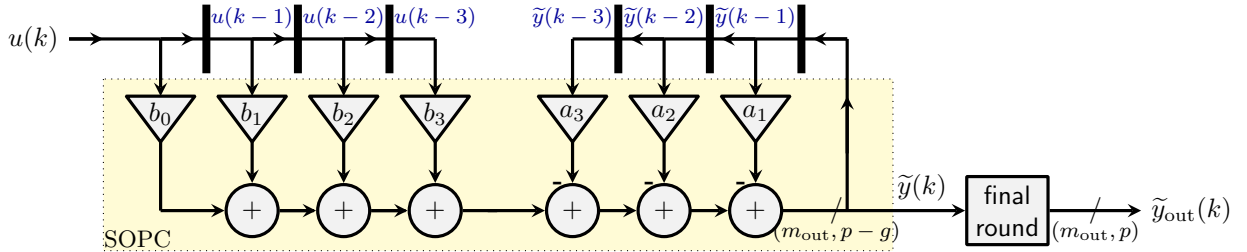


Figure 1: Abstract architecture for the direct form realization of an LTI filter

#### 1.3.2 State-space representation: a recurrence to define an infinite response

This type of realization consists in expressing the evolution of a system considering its state at time  $k$ . In continuous time, it is described by differential equations at first order. In discrete time (in which we are interested in), it is described by a simple recurrence:

$$\begin{cases} \mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k+1) = \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases} \quad (1.3.1)$$

Where  $\mathbf{x}(k) \in \mathbb{R}^{n_x}$  is the state vector,  $\mathbf{u}(k) \in \mathbb{R}^{n_u}$  is the input vector and  $\mathbf{y}(k) \in \mathbb{R}^{n_y}$  is the output vector, at time  $k$ . The matrices  $\mathbf{A} \in \mathbb{R}^{n_x \times n_x}$ ,  $\mathbf{B} \in \mathbb{R}^{n_x \times n_u}$ ,  $\mathbf{C} \in \mathbb{R}^{n_y \times n_x}$ , and  $\mathbf{D} \in \mathbb{R}^{n_y \times n_u}$ , with  $\mathbf{x}(0)$  are sufficient to describe an LTI filter, with convention  $\mathbf{x}(k) = \mathbf{u}(k) = 0$ ,  $\forall k < 0$ .

### 1.4 The Specialized Implicit Form (SIF): a unified representation

#### 1.4.1 Definition

The classical state-space representation is intuitive, but it doesn't takes into account the reality of implementation. The specialized implicit form (SIF) was introduced in [8] and is well detailed in Hilaire's PhD and

associated papers [5, 7], and is a good answer to this problem. Indeed, dealing with finite precision and error amplification, the order in which operations are done becomes crucial. Another motivation is to have a unique representation for any realization of LTI filters, that allows to compute every degradation measures instead of redevelopping them for each new realization. This form distinguishes computations done at one time from computations done the other times. As well as in a state-space,  $\mathbf{x}_i$  are state variables, but in addition to that,  $\mathbf{t}_i$  are intermediate variables. The SIF is described as following:

$$\begin{pmatrix} \mathbf{J} & \mathbf{0} & \mathbf{0} \\ -\mathbf{K} & \mathbf{I}_{n_x} & \mathbf{0} \\ -\mathbf{L} & \mathbf{0} & \mathbf{I}_{n_y} \end{pmatrix} \begin{pmatrix} \mathbf{t}(k+1) \\ \mathbf{x}(k+1) \\ \mathbf{y}(k) \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{M} & \mathbf{N} \\ \mathbf{0} & \mathbf{P} & \mathbf{Q} \\ \mathbf{0} & \mathbf{R} & \mathbf{S} \end{pmatrix} \begin{pmatrix} \mathbf{t}(k) \\ \mathbf{x}(k) \\ \mathbf{u}(k) \end{pmatrix} \quad (1.4.1)$$

With  $n_t$ ,  $n_x$ ,  $n_y$  and  $n_u$  the sizes of  $\mathbf{t}$ ,  $\mathbf{x}$ ,  $\mathbf{y}$  and  $\mathbf{u}$ , respectively.  $\mathbf{J}$ , is a lower triangular matrix, with diagonal entries equal to 1. The previous matrices have the following dimensions:

$$\begin{aligned} \mathbf{J} &\in \mathbb{R}^{n_t \times n_t}, \mathbf{M} \in \mathbb{R}^{n_t \times n_x}, \mathbf{N} \in \mathbb{R}^{n_t \times n_u}, \\ \mathbf{K} &\in \mathbb{R}^{n_x \times n_t}, \mathbf{P} \in \mathbb{R}^{n_x \times n_x}, \mathbf{Q} \in \mathbb{R}^{n_x \times n_u}, \end{aligned} \quad (1.4.2)$$

$$\mathbf{L} \in \mathbb{R}^{n_y \times n_t}, \mathbf{R} \in \mathbb{R}^{n_y \times n_x}, \mathbf{S} \in \mathbb{R}^{n_y \times n_u}, \quad (1.4.3)$$

The best way to understand the SIF may be to see it as an algorithm, each line of the equation 1.4.1 corresponding to a sequential step of the computation. The algorithm results as follows:

```

for int  $i = 0$  ;  $i \leq n_t$ ;  $i++$  do
  |  $\mathbf{t}_i(k+1) \leftarrow -\sum_{j < i} \mathbf{J}_{ij} \mathbf{t}_j(k+1) + \sum_{j=1}^{n_x} \mathbf{M}_{ij} \mathbf{x}_j(k) + \sum_{j=1}^{n_u} \mathbf{N}_{ij} \mathbf{u}_j(k)$ 
end
for int  $i = 0$  ;  $i \leq n_x$ ;  $i++$  do
  |  $\mathbf{x}_i(k+1) \leftarrow -\sum_{j=1}^{n_t} \mathbf{K}_{ij} \mathbf{t}_j(k+1) + \sum_{j=1}^{n_x} \mathbf{P}_{ij} \mathbf{x}_j(k) + \sum_{j=1}^{n_u} \mathbf{Q}_{ij} \mathbf{u}_j(k)$ 
end
for int  $i = 0$  ;  $i \leq n_y$ ;  $i++$  do
  |  $\mathbf{y}_i(k+1) \leftarrow -\sum_{j=1}^{n_t} \mathbf{L}_{ij} \mathbf{t}_j(k+1) + \sum_{j=1}^{n_x} \mathbf{R}_{ij} \mathbf{x}_j(k) + \sum_{j=1}^{n_u} \mathbf{S}_{ij} \mathbf{u}_j(k)$ 
end

```

**Algorithm 1:** Computation of SIF outputs from inputs

Here, it is important to see that the form of  $\mathbf{J}$  allows to compute the  $\mathbf{t}_i$  sequentially. The algorithm can then be described as follows:

```

for int  $i = 0$  ;  $i \leq n_t$ ;  $i++$  do
  |  $\mathbf{t}_i(k+1) \leftarrow -\mathbf{J}'_i \mathbf{t}(k+1) + \mathbf{M}_i \mathbf{x}(k) + \mathbf{N}_i \mathbf{u}(k)$ 
end
for int  $i = 0$  ;  $i \leq n_x$ ;  $i++$  do
  |  $\mathbf{x}_i(k+1) \leftarrow -\mathbf{K}_i \mathbf{t}(k+1) + \mathbf{P}_i \mathbf{x}(k) + \mathbf{Q}_i \mathbf{u}(k)$ 
end
for int  $i = 0$  ;  $i \leq n_y$ ;  $i++$  do
  |  $\mathbf{y}_i(k+1) \leftarrow -\mathbf{L}_i \mathbf{t}(k+1) + \mathbf{R}_i \mathbf{x}(k) + \mathbf{S}_i \mathbf{u}(k)$ 
end

```

**Algorithm 2:** Simplified matricial algorithm

With  $\mathbf{J}' = \mathbf{J} - \mathbf{I}_{n_t}$ .

Values of the vector  $\mathbf{t}(k+1)$  are computed and used at the same iterations, so they are not kept in memory. As the equation 1.4.1 is mostly full of zeros, it is more convenient to use it's compressed formulation, wich is denoted as the  $\mathbf{Z}$ :

$$\mathbf{Z} = \begin{pmatrix} -\mathbf{J} & \mathbf{M} & \mathbf{N} \\ \mathbf{K} & \mathbf{P} & \mathbf{Q} \\ \mathbf{L} & \mathbf{R} & \mathbf{S} \end{pmatrix} \quad (1.4.4)$$

The community usually takes  $-\mathbf{J}$  for simplicity within further computations.

The SIF can of course be transformed into an equivalent ABCD (classic state-space) form, which gives:

$$\begin{aligned} \mathbf{A}_Z &= \mathbf{K}\mathbf{J}^{-1}\mathbf{M} + \mathbf{P}, & \mathbf{B}_Z &= \mathbf{K}\mathbf{J}^{-1}\mathbf{N} + \mathbf{Q}, \\ \mathbf{C}_Z &= \mathbf{L}\mathbf{J}^{-1}\mathbf{M} + \mathbf{R}, & \mathbf{D}_Z &= \mathbf{L}\mathbf{J}^{-1}\mathbf{N} + \mathbf{S}, \end{aligned} \quad (1.4.5)$$

with:

$$\begin{aligned} \mathbf{A}_Z &\in \mathbb{R}^{n_x \times n_x}, \mathbf{B}_Z \in \mathbb{R}^{n_x \times n_u}, \\ \mathbf{C}_Z &\in \mathbb{R}^{n_y \times n_x}, \mathbf{D}_Z \in \mathbb{R}^{n_y \times n_u}, \end{aligned} \quad (1.4.6)$$

#### 1.4.2 Workflow for filter implementation

The choice of the SIF realization is not concerned by this work. This choice is up to the user, who may have a lot of good reasons not to design a realization that seem intuitive for us. The choice of the realization is a job done at LIP6 in the PEQUAN team, under the metalibm ANR project. The present work is just a hardware backend for the optimization part. Of course, optimizations are available at every level. The expected use of this work is shown on figure 2.

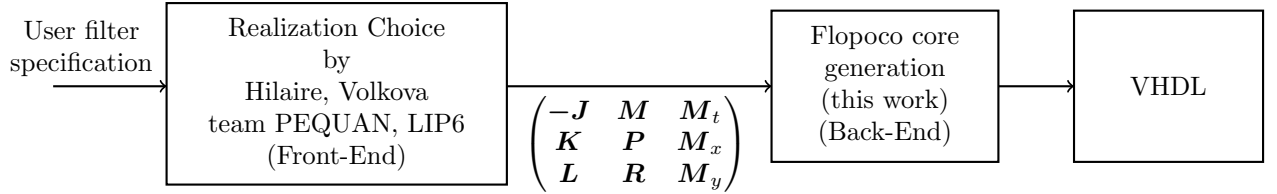


Figure 2: Workflow overview of tools usage

## 2 Architecture generation

### 2.1 Implementation: what can be done in practise?

#### 2.1.1 The FloPoCo Framework: computing just right

FloPoCo is a C++ framework [4] which first purpose is to generate floating point cores in VHDL. It targets the configuration of FPGAs as computing units, and could also be useful dealing with ASIC specification. Indeed, FPGAs are well used for embedded computing, because they offer full reconfigurability and a relatively good balance between computation power and price.

More information is available at <http://flopoco.gforge.inria.fr/>

This work proposes to implements SIFs using the SOPCs arithmetical core from FloPoCo. SOPC stands for Sum of Products by Constants, and is a KCM-multiplier based architecture. The details of the SOPC architecture has been described in [3].

#### 2.1.2 KCM multipliers: an architecture to efficiently compute products on FPGAs

To be quick, a KCM multiplier is designed to fully use the power of look-up tables (LUTs) in an FPGA. To do so, it partitions the constant into chunks which match the size of a LUT. Then, tabulating the multiplications ends up with a final sum to compute the product. This has been first described by K.Chapman in [1]. The detail of implementation is visible on figure 3

#### 2.1.3 SOPCs: exploiting full parallel architectures

The SOPC architecture keeps the same idea. The difference resides in the mutualization of the summation architecture through the  $n_c$  products. This summation is performed within a bit-heap based architecture implemented by flopoco. The detail is visible on figure 4.

To generate such cores, FloPoCo needs a number format specification, that is, the most significant bit (msb) and last significant bit (lsb) for the inputs and output of each operator to generate.

#### 2.1.4 Precision concerns

In SOPCs architectures, the accuracy is deducible from the inputs/outputs specifications and the size of the constants. This is described in [3].

Dealing with feedback inputs to build an SOPC-based filter, the question of precision is more complicated. Indeed, when outputs loop back on inputs, as soon as the problem is in finite precision, the error is amplified by a certain amount, depending on the coefficients, at each pass through the filter.

In this case the solution in the industry is to build an equivalent FIR filter by resolving the state-space recurrence. This leads to build an accurate hardware, but at the price of a huge waste of logic. The present work tries to answer this problem, trying to compute just right, keeping the recurrence and saving hardware.

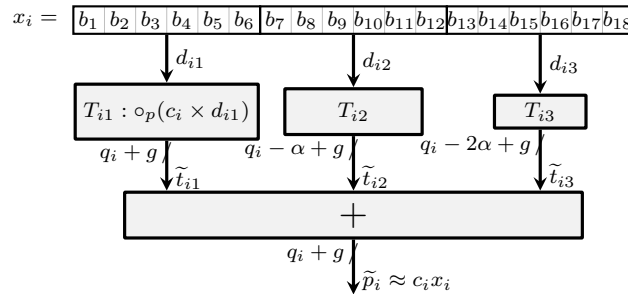
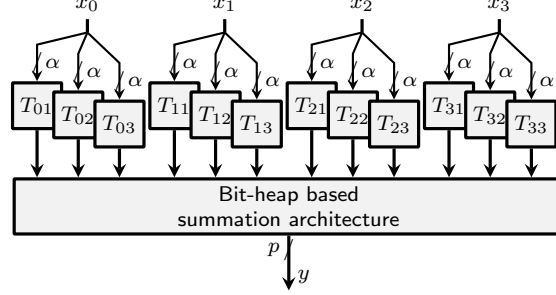


Figure 3: The FixRealKCM method when  $x_i$  is split in 3 chunks




 Figure 4: KCM-based SPOC architecture for  $n_c = 4$ , each input being split into 3 chunks

The main idea to dimension filters is to consider the total error as a single filter. The result of this filter is then added to the “perfect” filter to get the final output.

The two incoming parts are calculations from Lopez’s PhD [6], with adaptations in the architecture context. Indeed, the original calculations are done for fixed precision, because they are expected to be done in software, with fixed word sizes. The challenge here is to produce architecture with arbitrary precision, that is, choosing the best size for each basic operation.

### 2.1.5 Size computation in finite precision

When there are many potential realizations for a single LTI filter, the choice of one realization among the others is very related to the error analysis in output. The precision of our computations can then be defined so that the following condition is satisfied:

$$\varepsilon_{y_i} < 2^{-lsb_{y_i}} \quad \forall i \in [1, n_y] \quad (2.1.1)$$

With  $lsb_{y_i}$  the last significant bit of the output  $y_i$ , and  $\varepsilon_{y_i}$  the error on the computation of the output  $y_i$

Following the same model, an error term for each SOPC is involved in the computation. All these error terms will be functions of the most significant and last significant bits (respectively msb and lsb) of every input, output and intermediate signal. So there is a need for computing every  $msb_{t_i}$ ,  $lsb_{t_i}$ ,  $msb_{x_i}$ ,  $lsb_{x_i}$  and  $msb_{y_i}$ ,  $lsb_{y_i}$ ,

This leads to the definition of errors introduced by the SOPCs architectures, based on the SIF computation algorithm:

$$\varepsilon_t(k) = -J't^*(k+1) + Mx^*(k) + Nu(k) - SOPC(J', M, N, t^*(k), x^*(k), u(k)), \quad (2.1.2)$$

$$x^*(k+1) = Kt^*(k+1) + Px^*(k) + Qu(k) - SOPC(K, P, Q, t^*(k), x^*(k), u(k)), \quad (2.1.3)$$

$$y^*(k+1) = Lt^*(k+1) + Rx^*(k) + Su(k) - SOPC(L, R, S, t^*(k), x^*(k), u(k)), \quad (2.1.4)$$

where  $SOPC(L, R, S, t^*(k), x^*(k), u(k))$  is the effective computation (with errors) of the sum of products.

Let's define the following vectors:

$$\boldsymbol{\varepsilon}(k) = \begin{pmatrix} \boldsymbol{\varepsilon}_t(k) \\ \boldsymbol{\varepsilon}_x(k) \\ \boldsymbol{\varepsilon}_y(k) \end{pmatrix} = \begin{pmatrix} \varepsilon_{t_1}(k) \\ \varepsilon_{t_2}(k) \\ \vdots \\ \varepsilon_{t_{n_t}}(k) \\ \varepsilon_{x_1}(k) \\ \varepsilon_{x_2}(k) \\ \vdots \\ \varepsilon_{x_{n_x}}(k) \\ \varepsilon_{y_1}(k) \\ \varepsilon_{y_2}(k) \\ \vdots \\ \varepsilon_{y_{n_y}}(k) \end{pmatrix} \quad (2.1.5)$$

### 2.1.6 Error analysis

Considering a filter  $\mathcal{H}$  and it's SIF, following the algorithm 2, in finite precision, leads to the following equations:

$$\mathbf{t}^*(k+1) = -\mathbf{J}'\mathbf{t}^*(k+1) + \mathbf{M}\mathbf{x}^*(k) + \mathbf{N}\mathbf{u}(k) + \boldsymbol{\varepsilon}_t(k) \quad (2.1.6)$$

$$\mathbf{x}^*(k+1) = \mathbf{K}\mathbf{t}^*(k+1) + \mathbf{P}\mathbf{x}^*(k) + \mathbf{Q}\mathbf{u}(k) + \boldsymbol{\varepsilon}_x(k) \quad (2.1.7)$$

$$\mathbf{y}^*(k+1) = \mathbf{L}\mathbf{t}^*(k+1) + \mathbf{R}\mathbf{x}^*(k) + \mathbf{S}\mathbf{u}(k) + \boldsymbol{\varepsilon}_y(k) \quad (2.1.8)$$

Here  $\mathbf{t}^*$ ,  $\mathbf{x}^*$  and  $\mathbf{y}^*$  are the computed vectors, so with computations errors.

As  $\boldsymbol{\varepsilon}$ -errors come only from the SOPCs cores, another error term will symbolize the total error. Let's denote:

$$\boldsymbol{\delta t}_i(k+1) = \mathbf{t}_i^*(k) - \mathbf{t}_i(k) \quad (2.1.9)$$

$$\boldsymbol{\delta x}_i(k+1) = \mathbf{x}_i^*(k) - \mathbf{x}_i(k) \quad (2.1.10)$$

$$\boldsymbol{\delta y}_i(k+1) = \mathbf{y}_i^*(k) - \mathbf{y}_i(k) \quad (2.1.11)$$

the total error at instant k, considering computations errors and loopback, for  $\mathbf{t}$ ,  $\mathbf{x}$  and  $\mathbf{y}$ . The equations, corresponding to an algorithm, become:

$$\boldsymbol{\delta t}(k+1) = -\mathbf{J}'\boldsymbol{\delta t}(k+1) + \mathbf{M}\boldsymbol{\delta x}(k) + \boldsymbol{\varepsilon}_t(k) \quad (2.1.12)$$

$$\boldsymbol{\delta x}(k+1) = \mathbf{K}\boldsymbol{\delta t}(k+1) + \mathbf{P}\boldsymbol{\delta x}(k) + \boldsymbol{\varepsilon}_x(k) \quad (2.1.13)$$

$$\boldsymbol{\delta y}(k+1) = \mathbf{L}\boldsymbol{\delta t}(k+1) + \mathbf{R}\boldsymbol{\delta x}(k) + \boldsymbol{\varepsilon}_y(k) \quad (2.1.14)$$

This new algorithm corresponds here to the algorithm of the SIF of a filter  $\mathcal{H}_{\boldsymbol{\varepsilon}}$ , which describes the behaviour of computation errors at time k on the output. The linearity condition allows to decompose the real  $\mathcal{H}^*$  filter in two distinct filters:

- $\mathcal{H}$  the absolute filter in infinite precision
- $\mathcal{H}_{\boldsymbol{\varepsilon}}$  the error filter

According to 1.4.4:

$$\mathbf{Z}_{\boldsymbol{\varepsilon}} = \begin{pmatrix} -\mathbf{J} & \mathbf{M} & \mathbf{M}_t \\ \mathbf{K} & \mathbf{P} & \mathbf{M}_x \\ \mathbf{L} & \mathbf{R} & \mathbf{M}_y \end{pmatrix} \quad (2.1.15)$$

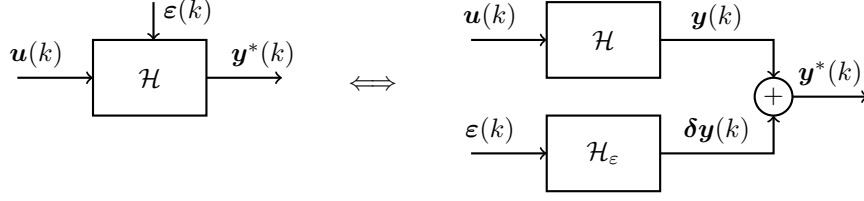


Figure 5: A signal view of the error propagation with respect to the ideal filter

with:

$$M_t = (\mathbf{I}_{n_t} \quad \mathbf{0}_{n_t \times n_x} \quad \mathbf{0}_{n_t \times n_y}), \quad (2.1.16)$$

$$M_x = (\mathbf{0}_{n_x \times n_t} \quad \mathbf{I}_{n_x} \quad \mathbf{0}_{n_x \times n_y}), \quad (2.1.17)$$

$$M_y = (\mathbf{0}_{n_y \times n_t} \quad \mathbf{0}_{n_y \times n_x} \quad \mathbf{I}_{n_y}), \quad (2.1.18)$$

$\mathcal{H}_\varepsilon$  is a filter with  $(n_t + n_x + n_u)$  inputs and  $n_y$  outputs.

*Proposition 1.* The transfert function of filter  $\mathcal{H}_\varepsilon$ , denoted  $\mathbf{H}_\varepsilon$ , is defined as follows:

$$\mathbf{H}_\varepsilon := \mathbf{C}_Z(z\mathbf{I}_n - \mathbf{A}_Z)^{-1}\mathbf{M}_1 + \mathbf{M}_2 \quad \forall z \in \mathbb{C} \quad (2.1.19)$$

with  $\mathbf{A}_Z$  and  $\mathbf{C}_Z$  the matrices defined by 1.4.5 and

$$\mathbf{M}_1 = (\mathbf{K}\mathbf{J}^{-1} \quad \mathbf{I}_{n_x} \quad \mathbf{0}), \quad \mathbf{M}_2 = (\mathbf{L}\mathbf{J}^{-1} \quad \mathbf{0} \quad \mathbf{I}_{n_y}), \quad (2.1.20)$$

The demonstration is well detailed in Lopez' PhD [6].

*Corollary 1.* Considering a filter  $\mathcal{H}$ ,  $\varepsilon(k)$  the vector of computation errors at time  $k$  in the finite precision of  $\mathcal{H}$ , and  $\mathcal{H}_\varepsilon$  the error filter associated to  $\mathcal{H}$ . The behaviour of error can be described from  $\varepsilon(k)$  and  $\mathcal{H}_\varepsilon$ . The error is considered as an interval vector, denoted by its center and radius  $\langle \varepsilon_m, \varepsilon_r \rangle$  and the interval vector of global error  $\delta \mathbf{y}$ , denoted  $\langle \delta \mathbf{y}_m, \delta \mathbf{y}_r \rangle$ .

In practise, all inputs in our case are centered around zero, which is not the case in the control community, where this notation makes sense. So,  $\varepsilon_m = 0$ . The results are the following:

$$\delta \mathbf{y}_m = 0 \quad (2.1.21)$$

$$\delta \mathbf{y}_r = \langle \langle \mathcal{H}_\varepsilon \rangle \rangle_{wcpq} \cdot \varepsilon_r \quad (2.1.22)$$

In the following  $\delta \mathbf{y}_m$  won't be taken into account.

Let's define:

$$n' = n_t + n_x + n_y$$

and:

$$\mathbf{v}' = \begin{pmatrix} \mathbf{t}(k+1) \\ \mathbf{x}(k+1) \\ \mathbf{y}(k) \end{pmatrix} \quad (2.1.23)$$

Then, following Lopez's computations, precisions are derived for every intermediate step:

$$|\delta \mathbf{y}_i| \leq \sum_{j=1}^{n'} |\langle \langle \mathcal{H}_\varepsilon \rangle \rangle_{i,j}| \cdot 2^{lsb_{v'_j}} \quad (2.1.24)$$

To formalize with a matricial formulation:

$$|\delta \mathbf{y}| \leq |\langle \langle \mathcal{H}_\varepsilon \rangle \rangle| \cdot 2^{lsb_{v'}} \quad (2.1.25)$$

To satisfy the condition 2.1.1,  $\xi$  is defined as the minimal error the user wants, which gives:

$$|\langle\langle\mathcal{H}_\epsilon\rangle\rangle| \cdot 2^{lsb_{v'}} < \xi \quad (2.1.26)$$

That is,

$$\mathfrak{A} \cdot 2^{lsb_{v'} - msb_{v'} - 1} < 1_{n_y} \quad (2.1.27)$$

Where:

$$\mathfrak{A}_{i,j} = |\langle\langle\mathcal{H}_\epsilon\rangle\rangle_{i,j}| \cdot \frac{2^{msb_{v'_j} + 1}}{\xi_i} \quad (2.1.28)$$

So, all  $lsb_{v_i}$  from all  $msb_{v_i}$  can be deduced from .  $msb_{v_i}$  can be deduced directly from  $\langle\langle\mathcal{H}\rangle\rangle_{wcpq}$ , choosing and computing each operator precision is possible.

### 2.1.7 Worst-case peak gain computation

Computing the worst-case peak gain accurately in finite precision is a non-trivial problem. A state of the art implementation is being developped at LIP6 by Anastatsi Lozanova Volkova from Hilaire's team (PEQUAN) [9]. This work will use it as soon as it is available.

## 2.2 Architecture generation algorithm

```
computePrecisions([msbs,lsbs][]) //get the matrix of msbs lsbs, functions of the wcpq for i=1;
i=Z.size(); i++ do
    row[ ]= Z[i][ ] //pick first row of Z
    for j=1; j=1; j=Z.size() j++ do
        | assign(SOPC[i], row[j], "T","X","U",[msbs,lsbs][i][j])
    end
    Second pass for wiring.
end
```

An example of implementation for a real-life case is given on figure 6.

On the figure 6, the colors in the matrix represent the different steps of computaions:

- blue for  $t$  computations
- green for  $x$  computations
- green for  $y$  computations

On the architecture scheme:

- grey background for x registers
- purple for loopback from the x registers
- orange for inputs

A small remark is needed here: most of coefficients of the  $\mathbf{J}$  coefficient are null in the original case specification. They are just kept under this form to show how the algorithm works and how the architecture is built.

## 2.3 Particular Forms, degenerated cases

### 2.3.1 ABCD Form: still achievable

The ABCD Form can be considered as a degenerated form of the SIF, with  $n_t = 0$ . The algorithm will work in this case too.

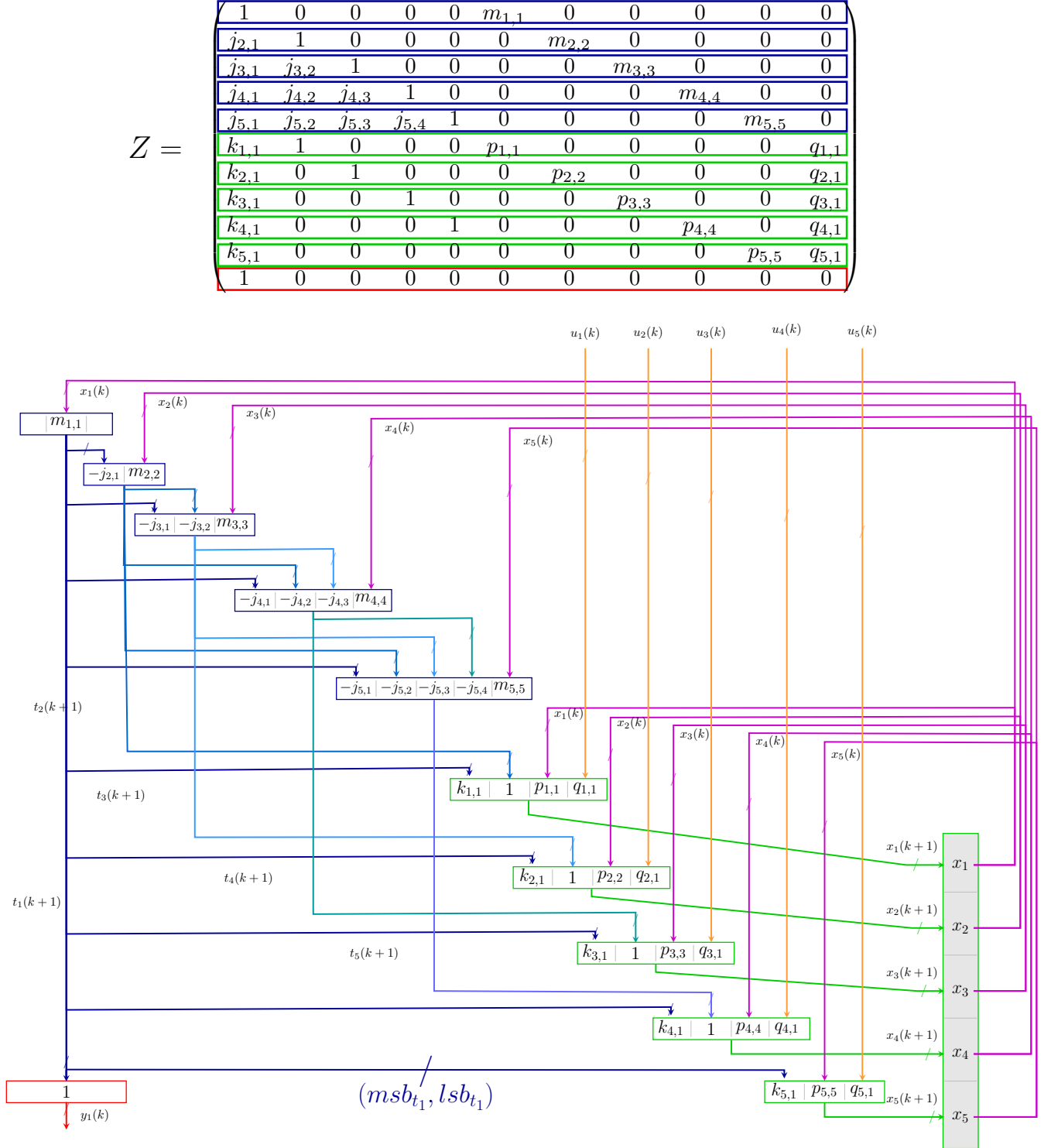


Figure 6: Architecture generation for implementing a SIF (example from a rho-DFII filter), with  $n_t = 5$ ,  $n_x = 5$  and  $n_y = 1$

### 2.3.2 When $n_x = 0$

When  $n_x = 0$ , the interest of using a SIF is of course very limited. This case is equivalent to building a FIR, and no loopback is needed. Still, the algorithm will work, allocating only SOPCs operators.

## 2.4 Optimizations to do on this technology

### 2.4.1 Sparse matrices

The Z matrix of a SIF is most of the time sparse. So it is useful to remove zeros coefficients before allocating SOPCs. Indeed, it prevents useless inputs to be declared and can save a lot of hardware, although the HDL compiler might be able to optimize the hardware and remove “dead code”. Anyway, it is healthy to keep a low compile time (either in flopoco or in the HDL compiler). Keeping the VHDL clean is always important, first for debugging issues, but also for comprehensiveness. This is already done in the current implementation.

### 2.4.2 Power of two coefficients

Coefficients equal to one or a power of two in the Z matrix can be interpreted as simple wires instead of multiplications in the SOPC. Entries equal to powers of two in SOPCs may be replaced by simple additions with the result of the SOPC. Here, the question of what solution is the best in terms of hardware consumption should be investigated (speed is not concerned here because the speed is determined by the length of the loop).

## 2.5 Implementation details

This implementation work has been done in the FixFilter section of the FloPoCo project. It is for now about 850 lines and contains:

- the implentation of the SIF algorithm
- the specification parser
- the testing framework (emulate method)

The testing framework consists on pre-computing expected results using the sollya lib [2]. Then the FloPoCo framework generates a VHDL testbench that embed the hardware architecture and the test of expected results. As an exhaustive testbench is impossible to generate ( $2^{msbIn-lsbIn} 2^{msbIn-lsbIn}$  values to test for each input), FloPoCo then generates a user-defined number of random tests.

## 2.6 Filter specification interface

To communicate with our SIF operator in FloPoCo, simply pass the coefficients in command line as it is done for the direct form is not affordable. So, a simple file format has been defined, in collaboration with our co-workers at LIP6, to store all the coefficients. The format is defined as follows:

```
X l c
x_1_1 x_1_2 ... x_1_c
x_2_1 x_2_2 ... x_2_c
.
.
.
x_l_1 x_l_2 ... x_l_c
```

Where X is the name of the matrix ( $X \in \{ J, K, L, M, N, P, Q, R, S, T \}$ ), x.i.j is the coefficient (with  $i \in [1, l]$  and  $j \in [1, c]$ ), l is the number of lines and c the number of columns.

All the matrices are specified after each other in the file.

For now another file is used to specify sizes because we are still waiting for the wcpg-code.

---

## 3 Further work

### 3.1 Static analysis on coefficients

Detection of and removal of null coefficients is already implemented. Detecting power of two entries in the coefficients seem to be a quite simple improvement to implement. However, this part should be left to the implementation of SOPCs, and so KCM multipliers. Indeed, pushing this detection at the lowest level is easy and gives the opportunity to the lower operators to give a feedback about their actual accuracy. Detecting a power of two in a KCM multiplier with feedback on precision will help to build better calculations on msbs and lsbs.

### 3.2 Sub-filter detection

Detecting independent loops can be very interesting in the context of saving hardware. Considering a sub-filter with its own loop can permit to use its WCPG to compute the precisions for this sub-filter. This leads to another resolution with a precision specification that depends on the rest of the filter (logic after the sub-filter in the pipeline). However, this optimization can be done at the front-end step, which leads to an open question. Indeed, detecting such loops has issues both in front-end and in backend, and it is a non-trivial problem.

### 3.3 Precision calculations improvement

This work kept original calculations trying to adapt them to our context. However, lots of approximations are done through those calculations. Working on this part, going back over all the calculations could really improve the efficiency and size of all implementations.

### 3.4 File format re-specification

Specifying new formats could help to improve the visibility and the usability of this new operator. The idea is to avoid specification output formatting operations for the user. The implementation should recognize input matrices formatted in the Python or Matlab formats.

## Conclusion

This work tried to give an overview of LTI filters, and the considerations to take into account when trying to implement them. This report tried to adapt concepts and calculation from Lopez's and Hilaire's work to the context of generating architectures computing just right. This context is merging ideas from the communities of automatic, signal, and computer arithmetics, so a certain amount of time was spent on the unification of notations and conventions. Finally, the implementation of a parametric definition in the FloPoCo framework was started. This was permitted through the definition of a specification language and the integration of external code from LIP6. Further work will take into account many improvements and optimizations, such as:

- static analysis on coefficients
- sub-filter detection
- bounds and precision computation improvements
- format re-specification

## References

- [1] K.D. Chapman. Fast integer multipliers fit in FPGAs (EDN 1993 design idea winner). *EDN magazine*, 39(10):80, May 1993.
- [2] S. Chevillard, M. Joldeş, and C. Lauter. Sollya: An environment for the development of numerical codes. In K. Fukuda, J. van der Hoeven, M. Joswig, and N. Takayama, editors, *Mathematical Software - ICMS 2010*, volume 6327 of *Lecture Notes in Computer Science*, pages 28–31, Heidelberg, Germany, September 2010. Springer.
- [3] Florent de Dinechin, Matei Istean, and Abdelbassat Massouri. Sum-of-product architectures computing just right. In *Application-Specific Systems, Architectures and Processors (ASAP)*. IEEE, 2014.
- [4] Florent de Dinechin and Bogdan Pasca. Designing custom arithmetic data paths with FloPoCo. *IEEE Design & Test of Computers*, 28(4):18–27, July 2011.
- [5] Thibaut Hilaire. *Analyse et synthèse de l'implémentation de lois de contrôle-commande en précision finie*. PhD thesis, Université de Nantes, 2006.
- [6] Benoit Lopez. *Implémentation optimale de filtres linéaires en arithmétique virgule fixe*. PhD thesis, Université Paris VI, 2014.
- [7] P.Chevel T.Hilaire and J.F.Whidbornz. A unifying framework for finite wordlength realizations. In *IEEE*, 2007.
- [8] P.Chevel T.Hilaire and Y.Trinquet. Implicit state-space representation: a unifying framework for FWL implementation of LTI systems. In *Proc. of the 16th IFAC World Congress.*, 2005.
- [9] A. Volkova, T. Hilaire, and C. Lauter. Reliable evaluation of the Worst-Case Peak Gain matrix in multiple precision. In *IEEE Symposium on Computer Arithmetic*, 2015.