

# Architecture synthesis for linear time-invariant filters

Antoine Martinet

June 10, 2015

# Contents

<b>Introduction</b>	<b>3</b>
<b>1 Motivations</b>	<b>3</b>
<b>2 Tools for expressing filters and signal processing</b>	<b>3</b>
2.0.1 Definition of a signal . . . . .	3
2.0.2 Linear Time Invariant Filters (LTI filters) . . . . .	3
2.0.3 Impulse response . . . . .	3
2.0.4 Worst-Case Peak Gain (WCPG) of a Filter . . . . .	4
2.1 FIR and IIR . . . . .	4
2.2 Realizations . . . . .	5
2.2.1 Direct and transposed forms . . . . .	5
2.2.2 State-space representation . . . . .	5
2.2.3 FloPoCo . . . . .	5
2.3 Specialized Implicit Form (SIF) . . . . .	5
2.4 Implementation . . . . .	7
2.5 Size computation in finite precision . . . . .	7
2.5.1 Error analysis . . . . .	8
2.5.2 A working solution to this problem . . . . .	10
<b>3 Dimensionment</b>	<b>11</b>
3.1 . . . . .	11
3.2 Algorithm . . . . .	11
3.3 Particular Forms . . . . .	11
3.4 ABCD Form . . . . .	11
3.5 When $n_x = 0$ . . . . .	13
3.6 Optimizations . . . . .	13
3.6.1 Sparse matrices . . . . .	13
3.6.2 One entries . . . . .	13
<b>Conclusion</b>	<b>13</b>
<b>Appendix: Remainders about signal processing</b>	<b>14</b>
3.7 Notations . . . . .	14
3.8 Remainders about signal processing . . . . .	14
3.8.1 Impulse response . . . . .	15

---

## Introduction

Filters are nowadays essential tools for designing responsive systems. As they appear in both signal processing and command communities, we can see that the use of filters is now spreading in many domains, related to them, such as radio signal coding and decoding, image and sound processing, and so on.

Although most of all the potential filters can be computed in software, some interesting applications in hardware could help to accelerate such computations, that might be very long in software in particular cases. Moreover, it allows to design and reconfigure specific parts of the hardware when working with FPGAs.

My work at the CITI lab in the SOCRATE team was to design a parametric architecture generator for LTI filters, which will be very useful to design a software-defined radio. This work has been integrated to FloPoCo tool, whose main goal is to generate architecture cores for computing just right.

In this report, I will first present reminders about signal processing. Then I will tell a little more about LTI filters and their software implementation in finite precision. Finally, I will present the main aspects of the hardware implementation of such filters.

## 1 Motivations

The main motivation of this work is to be able to compute efficiently LTI Filters in hardware. A more precise problem is the loopback problem, that is, when a filter uses its own outputs as inputs. Indeed, as the output is infinite, the problems of precision induce infinitely amplified errors, which cause the result to be unusable. In this case, the solution retained in the industry is for now to unroll such loops, building a very heavy hardware and wasting logic. The purpose of this work is to be able to build a hardware computing just right, avoiding logic waste.

## 2 Tools for expressing filters and signal processing

### 2.0.1 Definition of a signal

The state of the art comes mainly from Lopez's PhD [4] and Hilaire's PhD [3].

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

### 2.0.2 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  $\mathbf{y} = \mathcal{H}(\mathbf{u})$ , of size  $\dim(\mathbf{y}) = n_y$ . When  $n_u = n_y = 1$ , we speak about Single Input Single Output (SISO) filters. In other cases, we speak about Multiple Input Multiple Output (MIMO) filters.

**Definition 2.** (*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:*

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

### 2.0.3 Impulse response

**Definition 3.** (*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 can be described as a sum of Dirac impulsions:

$$u = \sum_{i \geq 0} u(i) \delta_i$$

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

$$\delta(k) = \begin{cases} 1 & \text{when } k = l \\ 0 & \text{else} \end{cases} \quad (2.0.1)$$

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, we have  $\mathbf{h} \in \mathbb{R}^{n_y \times n_u}$  as the impulse response of  $\mathcal{H}$ .  $\mathbf{h}_{i,j}$  is the response on the  $i$ th output to the Dirac implusion on the  $j$ -th input. The precedent equation becomes:

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

#### 2.0.4 Worst-Case Peak Gain (WCPG) of a Filter

**Definition 4.** (Worst-Case Peak Gain) The worst case peak gain is defined as the maximum amplification possible over all potential inputs through the filter.

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

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).

### 2.1 FIR and IIR

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

$$\exists n \in \mathbb{N} | \forall k \geq n, h(k) = 0 \quad (2.1.1)$$

The smallest  $n$  verifying 2.1.1 is referred as the order of the filter. So a  $n$ -order FIR can be described by the following equation:

$$y(k) = \sum_{i=0}^n b_i u(k - i) \quad (2.1.2)$$

An IIR will be described as following:

$$y(k) = \sum_{i=0}^n b_i u(k - i) - \sum_{i=0}^n a_i y(k - i) \quad (2.1.3)$$

Here one can observe that the output at time  $k$  depends also on all previous  $n$  outputs (loopback). One can also see 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 2.1.3 by resolving the recurrence relation:

$$h(k) = \begin{cases} 0 & \text{when } k < 0 \\ b_k - \sum_{l=1}^n a_l h(k - l) & \text{when } 0 \leq k \leq n \\ \sum_{l=1}^n a_l h(k - l) & \text{when } n < k \end{cases} \quad (2.1.4)$$

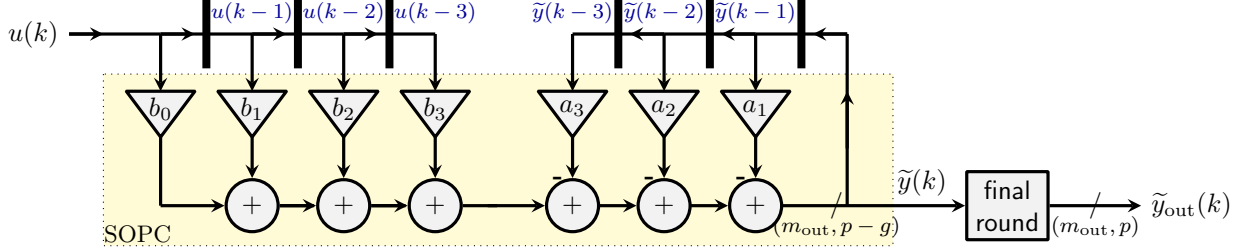


Figure 1: Abstract architecture for the direct form realization of an LTI filter

## 2.2 Realizations

**Definition 5.** (*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, order, 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, we will describe realizations as forms.

### 2.2.1 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 [4] [3]. The direct form has been implemented and in the FoPoCo project, you can see the hardware implementation an IIR on figure 2.2.1

### 2.2.2 State-space representation

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

$$\begin{cases} \mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k+1) = \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{cases} \quad (2.2.1)$$

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  $\mathbf{x}(k) = \mathbf{u}(k) = 0 \mid \forall k < 0$ .

### 2.2.3 FloPoCo

FloPoCo is a C++ framework which first purpose is to generate floating point cores in VHDL. Although it may have some pertinence in the ASIC design market, it main target is the configuration of FPGAs as computing units. Indeed, FPGAs are a rising market 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/>

## 2.3 Specialized Implicit Form (SIF)

he specialized implicit form (SIF) was introduced in [6] and is well detailed in [3, 5]. The classical state-space representation is intuitive, but it doesn't takes into account the reality of implementation. Indeed, dealing with finite precision and error amplification, the order in which operations are done becomes crucial. Another idea 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, we have  $\mathbf{x}_i$  as state variables, but in addition to that,  $\mathbf{t}_i$  are intermediate variables. The SIF is described as following:

$$\begin{pmatrix} \mathbf{J} & \mathbf{0} & \mathbf{0} \\ -\mathbf{K} & \mathbf{I}_{n_x} & \mathbf{0} \\ -\mathbf{L} & \mathbf{0} & \mathbf{I}_{n_y} \end{pmatrix} \begin{pmatrix} \mathbf{t}(k+1) \\ \mathbf{x}(k+1) \\ \mathbf{y}(k) \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{M} & \mathbf{N} \\ \mathbf{0} & \mathbf{P} & \mathbf{Q} \\ \mathbf{0} & \mathbf{R} & \mathbf{S} \end{pmatrix} \begin{pmatrix} \mathbf{t}(k) \\ \mathbf{x}(k) \\ \mathbf{u}(k) \end{pmatrix} \quad (2.3.1)$$

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. Then we have the following dimensions for the previous matrices:

$$\begin{aligned} \mathbf{J} &\in \mathbb{R}^{n_t \times n_t}, \mathbf{M} \in \mathbb{R}^{n_t \times n_x}, \mathbf{N} \in \mathbb{R}^{n_t \times n_u}, \\ \mathbf{K} &\in \mathbb{R}^{n_x \times n_t}, \mathbf{P} \in \mathbb{R}^{n_x \times n_x}, \mathbf{Q} \in \mathbb{R}^{n_x \times n_u}, \end{aligned} \quad (2.3.2)$$

$$\mathbf{L} \in \mathbb{R}^{n_y \times n_t}, \mathbf{R} \in \mathbb{R}^{n_y \times n_x}, \mathbf{S} \in \mathbb{R}^{n_y \times n_u}, \quad (2.3.3)$$

The best way to understand the SIF may be to see it as an algorithm, each line of the equation 2.3.1 corresponding to a sequential step of the computation. The algorithm results as follows:

```

for int  $i = 0 ; i \leq n_t ; i++$  do
    |  $\mathbf{t}_i(k+1) \leftarrow -\sum_{j=1}^{n_t} \mathbf{J}_{ij} \mathbf{t}_j(k+1) + \sum_{j=1}^{n_x} \mathbf{M}_{ij} \mathbf{x}_j(k) + \sum_{j=1}^{n_u} \mathbf{N}_{ij} \mathbf{u}_j(k)$ 
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 is mostly full of zeros, it is more convenient to use it's compressed formulation, wich is denoted as the  $\mathbf{Z}$ :

$$\mathbf{Z} = \begin{pmatrix} -\mathbf{J} & \mathbf{M} & \mathbf{N} \\ \mathbf{K} & \mathbf{P} & \mathbf{Q} \\ \mathbf{L} & \mathbf{R} & \mathbf{S} \end{pmatrix} \quad (2.3.4)$$

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

We can of course bind the SIF with the ABCD (classic state-space) form, which gives:

$$\begin{aligned} \mathbf{A}_Z &= \mathbf{K}\mathbf{J}^{-1}\mathbf{M} + \mathbf{P}, \quad \mathbf{B}_Z = \mathbf{K}\mathbf{J}^{-1}\mathbf{N} + \mathbf{Q}, \\ \mathbf{C}_Z &= \mathbf{L}\mathbf{J}^{-1}\mathbf{M} + \mathbf{R}, \quad \mathbf{D}_Z = \mathbf{L}\mathbf{J}^{-1}\mathbf{N} + \mathbf{S}, \end{aligned} \quad (2.3.5)$$

With:

$$\begin{aligned} \mathbf{A}_Z &\in \mathbb{R}^{n_x \times n_x}, \mathbf{B}_Z \in \mathbb{R}^{n_x \times n_u}, \\ \mathbf{C}_Z &\in \mathbb{R}^{n_y \times n_x}, \mathbf{D}_Z \in \mathbb{R}^{n_y \times n_u}, \end{aligned} \quad (2.3.6)$$

## 2.4 Implementation

In this work, we propose 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 [2]. 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, we end 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 2

The SOPC architecture keeps the same idea. The difference is, the summation architecture is mutualized through the  $n_c$  products in a bit-heap based architecture implemented by flopoco. The detail is visible on figure 3.

## 2.5 Size computation in finite precision

When there is many potential realizations for a single LTI filter, we can see that 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:

$$\varepsilon_{y_i} < 2^{-lsb_{y_i}} \quad \forall i \in [1, n_y] \quad (2.5.1)$$

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, we define an error quantum for each SOPC involved in the computation. All these error quantas will be functions of the most significant and last significant bits (respectively msb and lsb) of every input, output and intermediate signal.

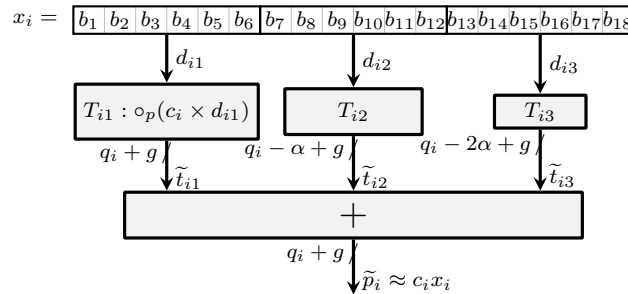


Figure 2: The FixRealKCM method when  $x_i$  is split in 3 chunks

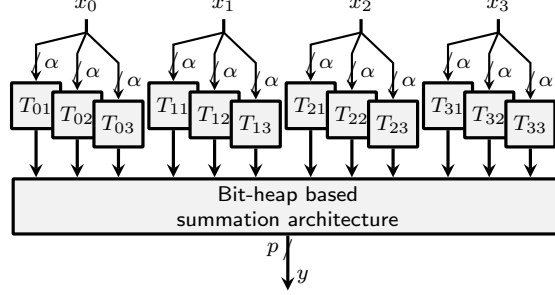


Figure 3: KCM-based SPOC architecture for  $n_c = 4$ , each input being split into 3 chunks

We define the following vectors:

$$\varepsilon(k) = \begin{pmatrix} \varepsilon_t(k) \\ \varepsilon_x(k) \\ \varepsilon_y(k) \end{pmatrix} = \begin{pmatrix} \varepsilon_{t_1}(k) \\ \varepsilon_{t_2}(k) \\ \vdots \\ \varepsilon_{t_{n_t}}(k) \\ \varepsilon_{x_1}(k) \\ \varepsilon_{x_2}(k) \\ \vdots \\ \varepsilon_{x_{n_x}}(k) \\ \varepsilon_{y_1}(k) \\ \varepsilon_{y_2}(k) \\ \vdots \\ \varepsilon_{y_{n_y}}(k) \end{pmatrix} \quad (2.5.2)$$

### 2.5.1 Error analysis

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

$$\mathbf{t}^*(k+1) = -\mathbf{J}'\mathbf{t}^*(k+1) + \mathbf{M}\mathbf{x}^*(k) + \mathbf{N}\mathbf{u}(k) + \varepsilon_t(k) \quad (2.5.3)$$

$$\mathbf{x}^*(k+1) = \mathbf{K}\mathbf{t}^*(k+1) + \mathbf{P}\mathbf{x}^*(k) + \mathbf{Q}_i\mathbf{u}(k) + \varepsilon_x(k) \quad (2.5.4)$$

$$\mathbf{y}^*(k+1) = \mathbf{L}\mathbf{t}^*(k+1) + \mathbf{R}\mathbf{x}^*(k) + \mathbf{S}_i\mathbf{u}(k) + \varepsilon_y(k) \quad (2.5.5)$$

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

As  $\varepsilon$ -errors come only from the SOPCs cores, we have another error quantum that will symbolize the total error. We denote

$$\delta\mathbf{t}(k+1) = \mathbf{t}_i^*(k) - \mathbf{t}_i(k) \quad (2.5.6)$$

$$\delta\mathbf{x}(k+1) = \mathbf{x}_i^*(k) - \mathbf{x}_i(k) \quad (2.5.7)$$

$$\delta\mathbf{y}(k+1) = \mathbf{y}_i^*(k) - \mathbf{y}_i(k) \quad (2.5.8)$$

the total error at instant  $k$ , considering computations errors and loopback, for  $\mathbf{t}$ ,  $\mathbf{x}$  and  $\mathbf{y}$ . We get then:

$$\delta\mathbf{t}(k+1) = -\mathbf{J}'\delta\mathbf{t}(k+1) + \mathbf{M}\delta\mathbf{x}(k) + \varepsilon_t(k) \quad (2.5.9)$$

$$\delta\mathbf{x}(k+1) = \mathbf{K}\delta\mathbf{t}(k+1) + \mathbf{P}\delta\mathbf{x}(k) + \varepsilon_x(k) \quad (2.5.10)$$

$$\delta\mathbf{y}(k+1) = \mathbf{L}\delta\mathbf{t}(k+1) + \mathbf{R}\delta\mathbf{x}(k) + \varepsilon_y(k) \quad (2.5.11)$$



This new algorithm corresponds here to the algorithm of the SIF of a filter  $\mathcal{H}_\varepsilon$ , 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
- $\mathcal{H}_\varepsilon$  the error filter

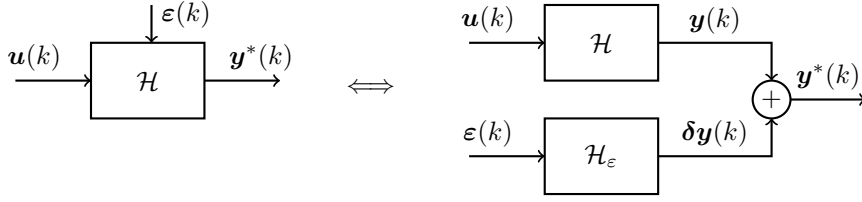


Figure 4: A signal view of the error propagation with respect to the ideal filter

According to 2.3.4, we have:

$$\mathbf{Z}_\varepsilon = \begin{pmatrix} -\mathbf{J} & \mathbf{M} & \mathbf{M}_t \\ \mathbf{K} & \mathbf{P} & \mathbf{M}_x \\ \mathbf{L} & \mathbf{R} & \mathbf{M}_y \end{pmatrix} \quad (2.5.12)$$

with:

$$\mathbf{M}_t = (\mathbf{I}_{n_t} \mathbf{0}_{n_t \times n_x} \mathbf{0}_{n_t \times n_y}), \quad (2.5.13)$$

$$\mathbf{M}_x = (\mathbf{0}_{n_x \times n_x} \mathbf{I}_{n_x} \mathbf{0}_{n_x \times n_y}), \quad (2.5.14)$$

$$\mathbf{M}_y = (\mathbf{0}_{n_y \times n_t} \mathbf{0}_{n_y \times n_x} \mathbf{I}_{n_y}), \quad (2.5.15)$$

$\mathcal{H}_\varepsilon$  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:

$$\mathbf{H}_\varepsilon : z \mapsto \mathbf{C}_Z (z \mathbf{I}_n - \mathbf{A}_Z)^{-1} \mathbf{M}_1 + \mathbf{M}_2 \quad \forall z \in \mathbb{C} \quad (2.5.16)$$

with  $\mathbf{A}_Z$  and  $\mathbf{C}_Z$  the matrices defined by 2.3.5 and

$$\mathbf{M}_1 = (\mathbf{K} \mathbf{J}^{-1} \quad \mathbf{I}_{n_x} \quad \mathbf{0}), \mathbf{M}_2 = (\mathbf{L} \mathbf{J}^{-1} \quad \mathbf{0} \quad \mathbf{I}_{n_y}), \quad (2.5.17)$$

The demonstration is well detailed in Lopez' phd.

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

$$\delta \mathbf{y}_m = \langle \langle \mathcal{H}_\varepsilon \rangle \rangle_{DC} \cdot \varepsilon_m \quad (2.5.18)$$

$$\delta \mathbf{y}_r = \langle \langle \mathcal{H}_\varepsilon \rangle \rangle_{wcp} \cdot \varepsilon_r \quad (2.5.19)$$

In practise, the interval arithmetic comes from the command community. In signal community, and moreover in computer arithmetics, we are used to compute on intervals centered around zero. Then, we can reasonably deduce that DC-gains are null, so the previous solution gives:

$$\delta \mathbf{y}_m = 0 \quad (2.5.20)$$

$$\delta \mathbf{y}_r = \langle \langle \mathcal{H}_\varepsilon \rangle \rangle_{wcp} \cdot \varepsilon_r \quad (2.5.21)$$

For the same reason (centered around zero), the interval becomes:

$$\delta \mathbf{y}_m = -\langle \langle \mathcal{H}_\epsilon \rangle \rangle_{wcp} \cdot \epsilon_r \quad (2.5.22)$$

$$\delta \mathbf{y}_r = \langle \langle \mathcal{H}_\epsilon \rangle \rangle_{wcp} \cdot \epsilon_r \quad (2.5.23)$$

Let's define:

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

and:

$$\mathbf{v}' = \begin{pmatrix} \mathbf{t}(k+1) \\ \mathbf{x}(k+1) \\ \mathbf{y}(k) \end{pmatrix} \quad (2.5.24)$$

Then, following Lopez's computations, we can derivate precisions for every intermediate step:

$$|\delta \mathbf{y}_i| \leq \sum_{j=1}^{n'} |\langle \langle \mathcal{H}_\epsilon \rangle \rangle_{i,j}| \cdot 2^{l_{v'_j}} \quad (2.5.25)$$

To formalize with a matricial formulation, we get:

$$|\delta \mathbf{y}| \leq |\langle \langle \mathcal{H}_\epsilon \rangle \rangle| \cdot 2^{lsb_{v'}} \quad (2.5.26)$$

To satisfy the condition 2.5.1, we can simply define  $\xi$  as the minimal error the user wants. We can then state:

$$|\langle \langle \mathcal{H}_\epsilon \rangle \rangle| \cdot 2^{lsb_{v'}} < \xi \quad (2.5.27)$$

That is,

$$\mathfrak{A} \cdot 2^{lsb_{v'} - msb_{v'} - 1} < \mathbf{1}_{n_y} \quad (2.5.28)$$

Where:

$$\mathfrak{A}_{i,j} = |\langle \langle \mathcal{H}_\epsilon \rangle \rangle_{i,j}| \cdot \frac{2^{msb_{v'_j} + 1}}{\xi_i} \quad (2.5.29)$$

$$(2.5.30)$$

**TODO:** check if this step is sufficient.

### 2.5.2 A working solution to this problem

A solution poposed in Lopez's phd [], is to consider the following reformulation of constraint 2.5.29. We can give a stronger majoration, as follows:

$$\frac{\mathfrak{A}_{i1}}{2^{msb_{v'_1} - lsb_{v'_1} + 1}} + \frac{\mathfrak{A}_{i2}}{2^{msb_{v'_2} - lsb_{v'_2} + 1}} + \dots + \frac{\mathfrak{A}_{in'}}{2^{msb_{v'_{n'}} - lsb_{v'_{n'}} + 1}} < 1, \quad \forall 1 \leq i \leq n_y \quad (2.5.31)$$

---

## 3 Dimensionment

### 3.1

In SOPCs architectures, the accuracy is deducible from the inputs/outputs specifications and the size of the constants.

This is described in

Dealing with feedback inputs, the question of precision is more complicated. Indeed, when results loop back to inputs, as soon as we are in finite precision, the error is amplified by a certain amount, depending on the coefficients, at each pass through the filter.

The main idea to dimension such 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.

In finite precision, sizes are constrained to be all the same. The demonstration of the size computation has been described in Lopez' PHD. The idea now is to see what we can do in arbitrary precision.

Here we have to compute each size at each step of the computation. Indeed, the WCPG is not useful for the first part of a FIR, as it has no loop. So we just need the WCPG for the second part, because it is just in this part of the circuit that there is a potential error amplification.

Direct and transposed forms are not directly transposable into SIF, but this problem is secondary.

### 3.2 Algorithm

```
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], TXU) //where TXU, is the indicator of the signal (this is just determined
        by the position of the coefficient)
    end
    Second pass for wiring.
end
```

An example of implementation for a real-life case is given on figure 5.

On the figure 5, the colors in the matrix represent the different steps of computations:

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

On the architecture scheme, we have:

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

A small precision is needed here. Here, most of coefficients of the  $\mathbf{J}$  coefficient are null. We just kept them under this form to show how the algorithm works and how the architecture is built.

### 3.3 Particular Forms

### 3.4 ABCD Form

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.

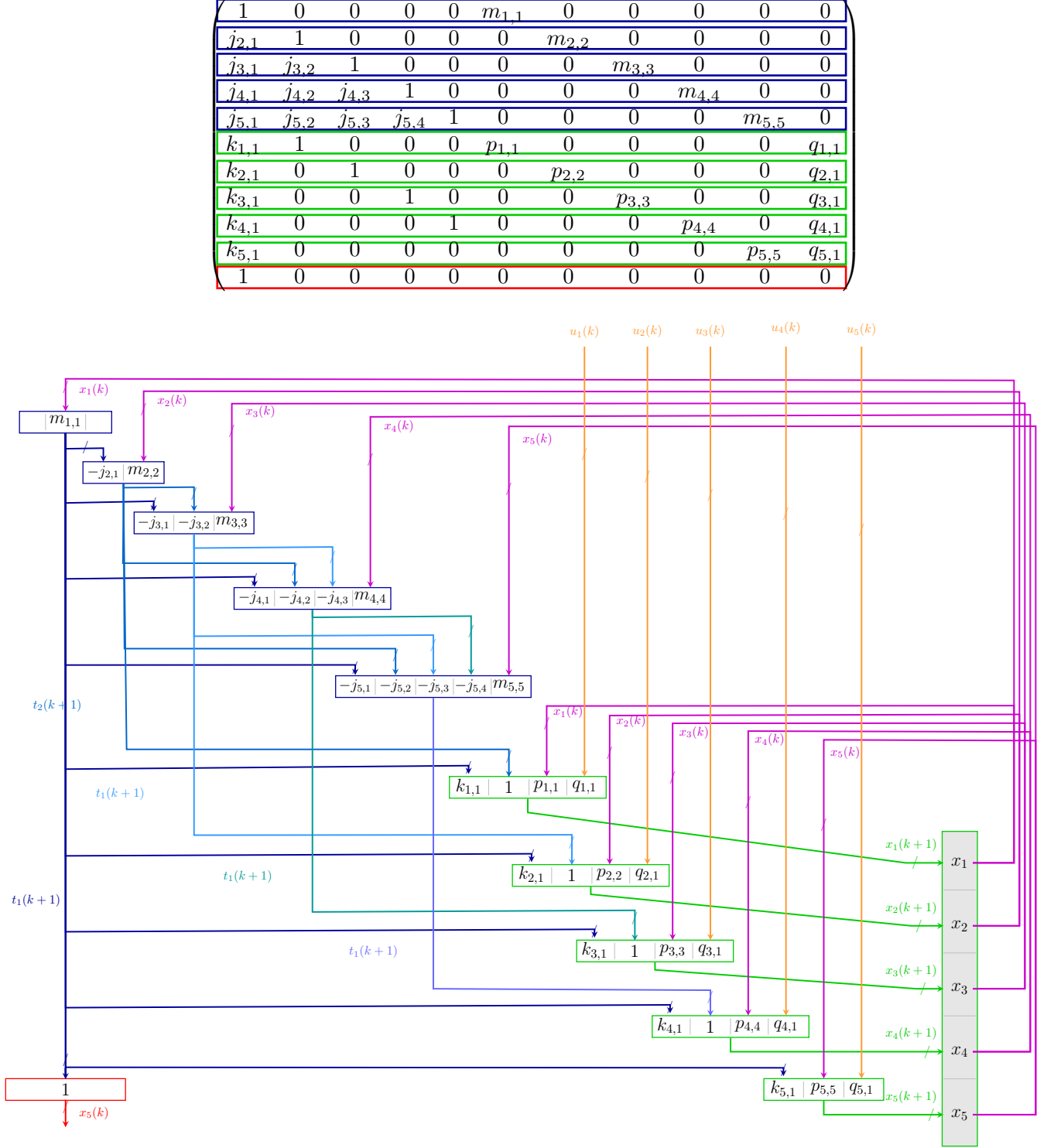


Figure 5: Architecture generation for implementing a SIF (example from a rho-DFII filter), with  $n_t = 5$ ,  $n_x = 5$  and  $n_y = 1$

### 3.5 When $n_x = 0$

When  $n_x = 0$ , the interest of using implicit form 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. The computation of precisions are then equivalent to the FIR precision computation, described in [].

## 3.6 Optimizations

### 3.6.1 Sparse matrices

The Z matrix of a SIF might be sparse in some degenerated cases. 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 more important, first for debugging issues, but also for comprehensiveness.

### 3.6.2 One entries

One entries in the Z matrix can be interpreted as simple wires instead of multiplications in the SOPC. So, we could eventually replace entries in SOPCs by simple additions with the result of the SOPC. Here, we should investigate to see what solution is the best in terms of hardware consumption (speed is not concerned here because the speed is determined by the length of the loop).

## Conclusion

## References

- [1] K.D. Chapman. Fast integer multipliers fit in FPGAs (EDN 1993 design idea winner). *EDN magazine*, 39(10):80, May 1993.
- [2] Abdelbassat Massouri Florent de Dinechin, Matei Istioan. Sum-of-product architectures computing just right. In *IEEE*, 2007.
- [3] Thibaut Hilaire. Analyse et synthèse de l'implémentation de lois de contrôle-commande en précision finie. In *PhD*, 2006.
- [4] Benoit Lopez. Implémentation optimale de filtre linéaire en arithmétique virgule fixe. In *PhD*, 2014.
- [5] P.Chevel T.Hilaire and J.F.Whidbornz. A unifying framework for finite wordlength realizations. In *IEEE*, 2007.
- [6] 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.

## Appendix: Reminders about signal processing

LTI filters in general are usually defined as sums of products. Several quantities are useful to understand their characteristics

### 3.7 Notations

During the entire report we will use several notations and conventions:

- in signal processing,  $t$  is the common notation for continue time and  $k$  the notation for discrete time. This report keeps this convention.
- $\langle\langle\mathcal{H}\rangle\rangle_{WCPG}$  is the worst case peak gain of the filter  $\mathcal{H}$
- $y$  is an output variable
- $x$  is a state variable
- $t$  is an intermediate variable
- 
- 

### 3.8 Remainders about signal processing

**Definition 6.** (*Signal*) Generally, a signal is a temporal variable, which takes a value from  $\mathbb{R}$  at each time  $t$ . We denote  $x(t)$  the value of the signal  $x$  at the instant  $t$ . When dealing with discrete time events, the time will be represented by  $k$ . Then we talk about  $x(k)$ , which is said to be a sample.  $\{x(k)\}_{k \geq 0}$  denotes all the values possible for the signal  $x$ . In the rest of this report, we will talk about vectors of signals  $\mathbf{x}$ , where  $\mathbf{x}(k) \in \mathbb{R}^n$

**Definition 7.** ( $\ell^1$  norm) The  $\ell^1$  norm of a scalar signal  $x$ , denoted  $\|x\|_{\ell^1}$ , is the sum of absolute values of  $x(k)$  at each instant  $k$ :

$$\|x\|_{\ell^1} = \sum_{k=0}^{+\infty} |x(k)| \quad (3.8.1)$$

This norm exists only if  $x$  is  $\ell^1$ -sommable, that is, if and only if the equation 3.8.1 converges.

**Definition 8.** ( $\ell^2$  norm) The  $\ell^2$  norm of a scalar signal  $x$ , denoted  $\|x\|_{\ell^2}$ , is defined as follows:

$$\|x\|_{\ell^2} = \sqrt{\sum_{k=0}^{+\infty} x(k)^2} \quad (3.8.2)$$

This norm exists only if  $x$  is square-sommable, that is, if and only if the equation 3.8.2 converges.

**Definition 9.** ( $\ell^\infty$  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:

$$\|x\|_{\ell^\infty} = \sup_{k \in \mathbb{N}} |x(k)| \quad (3.8.3)$$

### 3.8.1 Impulse response

**Definition 10.** (*Impulse Response*) A SISO filter may be defined by its impulse response, denoted  $h$ .  $h$  is the impulse response of  $\mathcal{H}$  to the impulsion of Dirac:

$$\delta(k) = \begin{cases} 1 & \text{when } k = 0 \\ 0 & \text{else} \end{cases} \quad (3.8.4)$$

Indeed each input can be described as a sum of Dirac impulsions:

$$u = \sum_{l \geq 0} u(l) \delta_l \quad (3.8.5)$$

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

$$\delta(k) = \begin{cases} 1 & \text{when } k = l \\ 0 & \text{else} \end{cases} \quad (3.8.6)$$

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:

$$y(k) = \sum_{l \geq 0} u(l) h(k - l) = \sum_{l=0}^k u(l) h(k - l) \quad (3.8.7)$$

This corresponds with the convolution product definition of  $u$  by  $h$ , denoted  $y = h * u$ .

Dealing with MIMO filters, we have  $\mathbf{h} \in \mathbb{R}^{n_y \times n_u}$  as the impulse response of  $\mathcal{H}$ .  $\mathbf{h}_{i,j}$  is the response on the  $i$ th output to the Dirac impulsion on the  $j$ th input.

The precedent equation becomes:

$$y_i(k) = \sum_{j=1}^{n_u} \sum_{l \geq 0}^k u_j(l) h_{i,j}(k - l) \quad \forall 1 \leq i \leq n_y \quad (3.8.8)$$