

*République Algérienne Démocratique et Populaire*  
*Ministère de L'Enseignement Supérieur et de la Recherche Scientifique*  
*Université Akli Mohand Oulhadj Bouira*



**FACULTE DES SCIENCES**  
**DEPARTEMENT GENIE ELECTRIQUE**

**MEMOIRE DE FIN D'ÉTUDES**  
pour l'obtention du diplôme de  
**MASTER**

Domaine : **Sciences et Techniques**  
Filière : **Génie Electrique**  
Spécialité :  **Systèmes électroniques complexes**

**Thème :**

---

*Développement d'une Solution au Niveau "Bit" au Problème de la  
Multiplication par une Constante*

---

*Soutenu publiquement le*

*Par :*

**Labandji Mourad**

**Promoteur:**

**M<sup>elle</sup> Sadji Aziza**  
**Mr Medjdoub Smail**

**Encadreur**

**Dr A. Kamel Oudjida**

**Jury:**

**- Dr. Madi Djemal MCA**  
**- Saoud Bilal**  
**- boucenna mohamed lamine**



*Je dédie ce modeste travail :*

*À ma très chère mère pour son amour et qui n'a jamais cessé de prier pour moi.*

*À mon très cher père, pour ses encouragements, son soutien, surtout pour son amour et son sacrifice afin que rien n'entrave le déroulement de mes études.*

*À Mes très chères sœurs.*

*À Mes très chers frères.*

*Toute ma famille,*

*À tous ceux qui contribuent de près ou de loin à la réalisation de ce travail,*

*À tous ceux qui sont proches de mon cœur.*

*L. MOURAD*

# Remerciements

*Nous remercions tout d'abord, le bon Dieu tout puissant de nous avoir donné la santé, la volonté, la force, le courage, et la puissance pour pouvoir surmonter les moments difficiles, et atteindre nos objectifs et sans lesquels notre projet n'aurait pas pu voir la lumière de ce jour.*

*J'exprime mes sincères remerciement à mon père, qui était mon soutien jusqu'au dernier jour de sa vie, merci bcp mon père, je remercie ma chère maman, elle est tous ce que j'ai maintenant dans ma vie.*

*Je souhaite vivement remercier et exprimer ma gratitude à mon encadreur Mr OUDJIDA pour m'avoir proposé ce sujet et qui je suis très reconnaissant pour ses remarques, ses orientations, son suivi pendant toute la réalisation de mon projet.*

*Je souhaite vivement remercier et exprimer ma gratitude à ma promotrice M<sup>elle</sup> SADJI qui j'exprime mon profonde reconnaissance pour les temps précieux qu'ils m'a consacré, ses suivis, ses nombreuses interventions qui ont permis l'aboutissement de ce travail et qui je suis très reconnaissant pour ses remarques, ses orientations, et ses conseils pendant toute la réalisation de mon projet et même pour la période d'étude pendant les trois années.*

*Je tiens également à remercier mon co-promoteur Mr MEJDOUB, de m'avoir orientés significativement dans mon travail.*

*Je dis mille fois merci, à ma famille, pour leur soutien et leurs encouragements ainsi qu'à tous ceux qui ont contribué de près ou de loin à la réalisation de ce modeste mémoire soit par un aide, un encouragement ou même un sourire.*

*Sans oublier de remercier très sincèrement mes chères amies pour tout quatre année d'étude en commun plein de bonheur, son encouragement énorme pour réussir a terminé ce projet.*

*Finalement, je remercie tout, les membres jury, devant qui j'ai l'honneur d'exposer mon travail, et qui ont pris la peine de lire avec soin ce mémoire pour juger son contenu.*

## *list des abréviations*

---

<b>ASIC</b>	: Application Specific Integrated Circuit
<b>Ath</b>	: Adder Depth, nombre maximum des opérations d'additions série de l'entrée à la sortie
<b>Avg</b>	: Average number of additions
<b>BASH</b>	: Bourne Again Shell
<b>BIGE</b>	: Bounded Inverse Graph Enumeration, algorithme SCM optimal
<b>CLB</b>	: Configurable Logic Bloc
<b>CSD</b>	: Canonical Signed Digit
<b>DAG</b>	: Directed Acyclic Graphs
<b>DBNS</b>	: Double Base Number System
<b>DSP</b>	: Digital-Signal-Processor/Processing
<b>FA</b>	: Full Adder
<b>FPGA</b>	: Field Programmable Gate Array
<b>FWL</b>	: Finite Word Length
<b>GCC</b>	: GNU Compiler Collection
<b>Hcub</b>	: Cumulative Benefit Heuristic
<b>HDL</b>	: Hardware Description Language
<b>IOB</b>	: Input Output Block
<b>LTI</b>	: Linear Time Invariant
<b>MAC</b>	: Multiply-And-Accumulate
<b>MCM</b>	: Multiple Constant Multiplication
<b>MEMS</b>	: Micro Electro Mechanical Systems
<b>NP-difficile</b>	: Non-deterministic Polynomial-time, problème complet pour la classe NP
<b>OM</b>	: Odd Multiple
<b>OS</b>	: Operating System
<b>PP</b>	: Partial Product
<b>RISC</b>	: Reduced Instruction Set Computer
<b>RNS</b>	: Residue Number System
<b>RTL</b>	: Register Transfer Level
<b>SCM</b>	: Single Constant Multiplication
<b>VHDL</b>	: Very high speed integrated circuit Hardware Description Language

## Résumé

Beaucoup d'applications dans le traitement de signal, d'image, des systèmes de contrôle ... etc. basés sur les systèmes LTI, ces systèmes sont généralement basés sur des multiplications par des constantes SCM/MCM qui consomme une quantité importante des ressources.

L'objectif de ce travail est le développement d'une solution au niveau bit pour la résolution de ces systèmes avec une optimisation des ressources consommées, la technique basée sur l'arithmétique RADIX-2r qui a des hautes capacités en matière de performance, et réduction du nombre des additions, de plus l'utilisation de cette heuristique au niveau bit permet plus de réduction sur les additionneurs de 1 bit, la solution décrivant une description VHDL robuste et optimisée.

La solution développée est une plateforme qui peut être utilisée par tout un système LTI, elle assure la construction de ce système avec une surface réduite, consommation de puissance réduite, et une vitesse élevée.

**Mots clé :** RADIX 2<sup>r</sup>, SCM, MCM, LTI, filtre RIF.

## Abstract

Many applications in signal processing, image processing, control systems, etc. based on LTI systems, these systems are generally based on multiplications by SCM / MCM constants that consume a significant amount of resources.

The objective of this work is the development of a solution at the bit level for the resolution of these systems with an optimization of the resources consumed, the technique based on arithmetic RADIX-2r which has high performance capacities, and reducing the number of additions, moreover the use of this heuristic at the bit level allows more reduction on the 1-bit adders, the solution describing a robust and optimized VHDL description.

The solution developed is a platform that can be used by an entire LTI system, it ensured the construction of this system with a reduced surface area, reduced power consumption, and high speed.

**Key words:** RADIX 2<sup>r</sup>, SCM, MCM, LTI, FIR filter.

## الملخص

العديد من التطبيقات في معالجة الإشارات ومعالجة الصور وأنظمة التحكم وما إلى ذلك. تعتمد غالباً على عمليات الضرب في ثابت أو مجموعة من الثوابت، التي تستهلك الكثير من الموارد.

والهدف من هذا العمل هو تطوير حل على مستوى البتات لهذه الأنظمة مع الاستفادة المثلى من الموارد المستهلكة، اعتماداً على تقنية RADIX 2<sup>r</sup> التي تتميز بالأداء العالي و الحد الأدنى من عمليات الجمع، إضافة إلى هذا فإن استخدامها على مستوى 1 بت تسمح بدرجة أكبر من الخفض، مع ملف وصفي VHDL قوي و أكثر دقة. هذا الحل هو منصة يمكن استعمالها في كل أنظمة LTI، فهو يضمن إنشاء نظام ذو مساحة منخفضة، سرعة في الأداء، مع استهلاك أقل للطاقة

الكلمات المفتاحية: RADIX 2<sup>r</sup>, SCM, MCM, LTI, FIR filter

# Table des figures

Figure I.1	format virgule fixe pour les nombres signés en complément à deux.....	7
Figure I.2	(a) Double et (b) simple précision.....	8
Figure I.3	Format RNS pour le nombre (8 7 5 3).....	13
Figure I.4	partitionnement de $(10599)_{10}$ dans $RADIX-2^3$ .....	15
Figure II.1	le nombre minimal d'additions pour $45 j \times X$ . ....	20
Figure II.2	Multiplication des constantes 81 et 23.....	21
Figure III.1	Partitionnement des segments d'une constante de n bits sous $RADIX-2^r$ .....	27
Figure III.2	Ordonnancement sequentiel de l'ensemble des produit partiels necessaire pour $RADIX-2^r$ .....	28
Figure III.3	Comparaison des limites supérieurs pour une constante de N bits.....	29
Figure III.4	Longueur de bits N d'une constante.....	30
Figure III.5	Le programme $RADIX-2^r$ .....	32
Figure III.6	Le recodage de 10599 et sa specification par l'algorithme $RADIX-2^r$ .....	33
Figure III.7	la solution SCM/MCM développée.....	35
Figure III.8	L'organigramme de la détection de la zone des multiples impaires.....	36
Figure III.9	Le principe de la structure chaine de liste.....	37
Figure III.10	L'organigramme de conversion et de stockage des coefficients.....	37
Figure III.11	architecture d'un exemple d'une solution SCM au niveau bloc d'additionneur.....	38
Figure III.12	L'organigramme de génération de la description VHDL.....	39
Figure III.13	Le processus de simulation de constante.....	40
Figure III.14	La description VHDL de $10599 \times X$ .....	40
Figure III.15	Le resultat de modelsim et le rapport de simulation de $1059 \times X$ .....	41
Figure IV.1	le principe d'optimisation au niveau d'additionneurs de 1 bit.....	42
Figure IV.2	la methode d'extension de sign appliqué dans $RADIX-2^r$ .....	43
Figure IV.3	les possibilité d'une addition/soustraction en $RADIX-2^r$ .....	46
Figure IV.4	l'organigramme de calcul du nombre total d'additionneurs 1 bit.....	47
Figure IV.5	organigramme de génération du code VHDL d'un SCM/MCM.....	48
Figure IV.6	un additionneur complet 1 bit.....	49
Figure IV.7	le circuit d'un additionneur soustracteur de 4 bits.....	50

Figure IV.8 L'architecture d'un SCM au niveau bit .....	51
Figure IV.9 l'achitecture d'un exemple SCM/MCM.....	51
Figure IV.10 test bench.....	52
Figure IV.11 principe de la verification fonctionnelle.....	52
Figure IV.12 quelques parties de script shell .....	53
Figure IV.13 le code test bench .....	55
Figure IV.14 compilation, élaboration et simulation du code VHDL .....	56
Figure IV.15 le script de simulation de modelsim.....	56
Figure IV.16 Comparaison entre le niveau bloc d'additionneurs et niveau bit.....	58
Figure IV.17 Variation d'un nombre d'additionneurs reduits en differents cas .....	58
Figure V.1 implementation d'un filtre RIF (a) forme direct (b) forme transposé.....	61
Figure V.2 la structure interne d'une FPGA .....	62
Figure V.3 L'architecture d'un filtre RIF .....	64
Figure V.4 l'organigramme de generation de description VHDL pour les filtres RIF .....	66
Figure V.5 les coefficients de filtres ainsi que leur recodage en RADIX-2 <sup>r</sup> .....	67
Figure V.6 Une partie de la description VHDL de filtre .....	68
Figure V.7 la generation des signaux d'entré pour la simulation.....	68
Figure V.8 les résultat de simulation de filtre.....	69

# List des tableaux

Tableau I.1 La table de vérité pour 10599 en RADIX-2r .....	14
Tableau I.2 nombre des chiffres requis pour chaque système pour la valeur (10599) <sub>10</sub> .....	16
Tableau I.3 Principales caractéristiques des systèmes numériques .....	16
Tableau III.1 Nombre moyen d'additions de RADIX-2r et CSD .....	29
Tableau III.2 Equations de RADIX-2r pour des constantes non négatives de différents longueur de bit.....	30
Tableau IV.1 la complexité de calcul des algorithmes SCM/MCM .....	44
Tableau IV.2 table de vérité d'un additionneur complet 1 bit.....	49
Tableau IV.3 table de vérité d'un inversement par XOR.....	50
Tableau IV.4 le nombre des additionneurs 1 bit consommés par les deux solutions .....	57
Tableau V.1 le traitement d'image dans une FPGA et un DSP.....	63

# Sommaire

<b>Introduction générale .....</b>	<b>1</b>
------------------------------------	----------

## **Chapitre I L'arithmétique binaire**

I.1 Introduction.....	3
I.2 L'arithmétique binaire .....	3
I.3 Les Formats de représentation des nombres .....	4
I.3.1 Format à virgule fixe .....	4
I.3.2 Format à virgule flottante .....	9
I.4 Les systèmes de représentation des nombres.....	10
I.4.1 La représentation numérique canonique des nombres signés (CSD).....	10
I.4.2 Système de numération à double Base (DBNS).....	11
I.4.3 Système de numération de résidus (RNS).....	11
I.4.4 Système de numération radix-2 <sup>r</sup> .....	13
I.4.5 Comparaison entre les systèmes de numération.....	15
I.5 Conclusion .....	17

## **Chapitre II Multiplication par une constante SCM/MCM**

II.1 Introduction .....	18
II.2 Les systèmes LTI.....	18
II.2.1 Le rôle des systèmes linéaires invariants .....	18
II.2.2 Formulation des systèmes LTI :.....	18
II.3 La contrainte de la multiplication par une constante .....	19
II.3.1 Multiplication de multiple-Constant (MCM).....	20
II.3.2 La formalisation du problème SCM .....	21
II.3.3 Les algorithmes SCM/MCM existants.....	22
II.3.4 Définition des mesures des algorithmes SCM/MCM .....	23
II.3.5 Les principales Limitations des algorithmes SCM/MCM existants : .....	23
II.4 Le Nouvel algorithme de recodage (RADIX-2 <sup>r</sup> ) .....	23
II.5 Conclusion .....	24

### **Chapitre III SCM/MCM au niveau bloc d'additionneur**

III.1 Introduction .....	25
III.2 RADIX-2 <sup>r</sup> pour la multiplication par une constante SCM/MCM .....	25
III.2.1 Nombre maximal d'additions pour une constante de N bits .....	27
III.2.2 Les caractéristiques de RADIX-2 <sup>r</sup> .....	31
III.3 Les spécifications d'entrée du système par l'application RADIX-2 <sup>r</sup> .....	32
III.4 representation de la solution developpée au niveau bloc d'additionneurs .....	33
III.4.1 Les langages utilisés .....	33
III.4.2 La description logicielle .....	34
III.4.3 Principe .....	35
III.4.4 Stratégie de translation des spécifications du système .....	35
III.4.5 Stratégie de génération de la description VHDL .....	38
III.5 Tests et simulations .....	39
III.6 Conclusion .....	41

### **Chapitre IV SCM/MCM au niveau bit**

IV.1 Introduction .....	42
IV.2 Le principe de RADIX-2 <sup>r</sup> SCM/MCM au niveau bit .....	42
IV.2.1 Extensions de signe .....	42
IV.2.2 Le nombre total d'additionneurs de 1 bit .....	43
IV.2.3 L'avantage de RADIX-2 <sup>r</sup> au niveau bit .....	44
IV.3 Le developpement d'une solution SCM/MCM au niveau bit .....	45
IV.3.1 Stratégie de calcul du nombre des additionneurs 1 bits .....	45
IV.3.2 La génération du code VHDL .....	47
IV.4 La description VHDL du SCM/MCM .....	49
IV.4.1 L'additionneur complet 1 bits (FA) .....	49
IV.4.2 L'additionneur soustracteur .....	49
IV.4.3 L'architecture d'un SCM/MCM .....	50
IV.5 La verification de fonctionnelle .....	51
IV.5.1 Le script shell .....	53

IV.5.2 Le Test bench .....	53
IV.5.3 Compilation, élaboration, et simulation .....	55
IV.5.4 La simulation par modelsim .....	56
IV.6 Resultats et discussion .....	57
IV.7 Conclusion .....	59

## **Chapitre V      Application du SCM/MCM aux filtres RIF**

V.1 Introduction .....	60
V.2 Les filtre RIF .....	60
V.3 Les plateformes de réalisation des filtres RIF .....	62
V.3.1 FPGA/ASIC .....	62
V.3.2 DSP .....	63
V.4 La structure d'un filtre RIF .....	64
V.4.1 Pipelining .....	64
V.5 La conception du programme générateur code VHDL pour les filtres RIF .....	65
V.5.1 Résultat et discussions .....	67
V.6 Conclusion .....	70
<b>Conclusion générale .....</b>	<b>71</b>

# Introduction générale

Le domaine des applications MEMS a la particularité d'être très sensible au bruit. Afin d'atteindre la précision souhaitée dans la préhension des objets micrométriques, le développement des méthodes avancées de contrôle est nécessaire. Ce dernier se traduit souvent par des lois de commande, des filtres et des algorithmes qui sont exécutés en temps réel sur les plates-formes qui sont volumineux et coûteux et présentent un grand inconvénient d'empêcher toute utilisation embarquée ou mobile (autonomie). Par conséquent, l'intégration matérielle de ces lois de commande et tous les problèmes y afférents constituent un véritable défi.

Notre travail aborde le problème d'optimisation des ressources matérielles au niveau bit de systèmes linéaires invariant dans le temps (Linear Time Invariant (LTI)) avec une longueur de mot finie (Finite Word Length (FWL)). L'effort d'optimisation de notre travail se situe au niveau algorithmique et non pas architecturale. Le plus grand défi est d'assurer un contrôle satisfaisant avec un nombre de ressources hardware le plus réduit possible. A ce stade deux optimisations distinctes sont possible bien que complémentaires : théorie de contrôle et l'arithmétique binaire. Seule cette dernière est associée à notre projet.

L'arithmétique binaire doit être à la fois assez rapide pour faire face à la dynamique rapide des MEMS, à faible consommation d'énergie pour des applications embarquées (utilisation d'une batterie), hautement évolutive (scalable) en terme de performances, et facilement prévisible pour donner une idée précise sur les ressources logiques nécessaires avant la mise en œuvre.

L'exploration d'un certain nombre d'arithmétiques binaires a montré que radix-2<sup>r</sup> est le meilleur candidat qui correspond aux exigences susmentionnées. Il a été pleinement exploité à concevoir des noyaux de multiplicateur ( $\times$ ) efficace, qui sont le véritable moteur des systèmes linéaires.

Pour mener bien notre projet, nous avons structuré notre mémoire en Cinq chapitres :

- Le premier chapitre inclut les fondements de l'arithmétique binaire. Tout d'abord les formats de nombre généralement utilisés format à virgule fixe et format à virgule flottante, leur précision et leur gamme dynamique respectives, Nous présentons ensuite les systèmes de numération les plus couramment utilisés. Nous insistons surtout sur les deux opérations arithmétique (+,  $\times$ ) requis par les systèmes linéaires.
- Le deuxième chapitre est consacré à l'opération de multiplication par une constante. Il présente une formulation des problèmes SCM/MCM (Single Constant Multiplication / Multiple Constant Multiplication) suivi par l'introduction d'une nouvelle heuristique appelée RADIX-2<sup>r</sup>.
- Le troisième chapitre traite le problème SCM/MCM au niveau block d'additionneurs (adder-bloc). Nous présentons notre solution développée dans le cadre de l'heuristique RADIX-2<sup>r</sup>. Nous montrons son efficacité par des résultats expérimentaux.

- Le quatrième chapitre traite par contre le problème SCM/MCM au niveau bit d'additionneur (full-adder). Nous analyserons notre solution par une étude aussi bien analytique qu'expérimentale.
- Dans le chapitre 5, nous appliquons les résultats de la recherche développée dans les chapitres précédents pour les filtres FIR .

A la fin on terminera par une conclusion générale et perspective.

# **Chapitre I :**

## **L'arithmétique binaire**

## I.1 Introduction

Ce chapitre aborde les fondements de l'arithmétique binaire. Nous décrivons les représentations les plus utilisés dans les systèmes numériques, leurs précisions, et leurs gammes dynamiques.

Ensuite, nous passerons vers les principaux systèmes de numération qui sont basés sur le format en virgule fixe, notamment le système de numération de résidus et le radix-2<sup>r</sup>, pour les opérations arithmétiques l'addition et de multiplication. Nous insisterons sur l'importance de ces deux opérations dans le calcul et la conception hardware des systèmes linéaires.

## I.2 L'arithmétique binaire

L'arithmétique binaire est un domaine mathématique basé principalement sur :

- ✓ L'étude des systèmes de représentation des nombres afin d'éliminer la redondance logique dans le recodage ;
- ✓ La recherche des algorithmes qui effectue efficacement les opérations arithmétiques ( $-$ ,  $+$ ,  $\times$ ,  $/$ ,  $\sqrt{\quad}$ ,  $a^n$  etc.) dans un système numérique donné ;
- ✓ L'exploration des meilleures techniques de mise en œuvre lorsque le système numérique est un DSP, un FPGA ou un ASIC.

En arithmétique binaire, les nombres sont représentés généralement à l'aide des deux symboles : 0 et 1.

Dans les systèmes de représentation des nombres, les symboles sont appelés chiffres, et le nombre de symboles est appelée "base" ou "radix".

Les nombres binaires sont exprimées en base 2, tandis que par exemple les nombres décimaux sont représentés en base 10, puisque dix chiffres sont utilisés (0, 1, ..., 9).

La notation utilisée couramment pour représenter des nombres est :  $(x)_y$ , où  $x$  est un nombre exprimé dans la base  $y$ .

**Exemple :**

$(10)_{10}$  représente le nombre dix dans le système décimal ;  $(10)_2$  représente le numéro deux dans le système binaire

*Remarque : les recodage octal (base 8), hexadécimal (base 16), etc. ..., ne servent qu'à faciliter la manipulation de longues chaînes de 0 et de 1. La mise en œuvre est réalisée en binaire (radix-2).*

Il existe différents types d'arithmétique binaire. Les caractéristiques de chaque calcul sont conditionnées par le format utilisé pour représenter les nombres. Les arithmétiques les plus souvent utilisées sont : arithmétique à virgule flottante et arithmétique à virgule fixe. Les implémentations et les optimisations menées dans chaque arithmétique sont différentes [1].

Le choix de l'arithmétique à utiliser est principalement dicté par la précision souhaitée et la gamme dynamique imposé par le problème à résoudre.

En outre, le niveau de complexité de l'arithmétique est déterminé par les opérations requises. Ce dernier dépend fondamentalement du type du système considéré : système linéaire ou non linéaire.

Un système linéaire est un modèle mathématique basé sur des opérations linéaires ( $-$ ,  $+$ ,  $\times$ ,  $/$ ). Il présente généralement des caractéristiques et des propriétés qui sont beaucoup plus simple et

plus facile à comprendre et à manipuler que le cas non linéaire qui nécessite des opérations complexes, tels que polynôme ou les fonctions trigonométriques.

### I.3 Les Formats de représentation des nombres

Il existe différents formats de représentation pour les nombres [4,8]. Les formats les plus utilisés sont résumés comme suit :

- **Format à virgule fixe** : propose une gamme limitée et/ou de précision, mais facile à implémenter. Très pratique pour les applications à haute vitesse et faible consommation de puissance. Il manipule des nombres entiers  $x \in I = \{-N, \dots, N\}$  ainsi que des nombres rationnels de la forme  $x = a/2^f$  ("binaire" rationnel),  $a \in I$  et  $f$  est un entier positif.
- **Format à virgule flottante** : il s'agit de l'approche la plus courante. Elle offre une gamme dynamique large avec une grande précision pour pouvoir manipuler des nombres extrêmement grands et petits, respectivement.

Toutefois, il est relativement difficile à implémenter. Il gère les nombres de la forme  $x \times b^E$ , où  $x$  est un nombre rationnel,  $b$  est la base en entier naturel, et  $E$  est un exposant entier naturel.

- **Format logarithmique** : représente les nombres par leurs signes et leurs logarithmes. Attrayant pour des applications qui ont besoin d'une faible précision et une plage dynamique étendue.
- **Format rationnel** : se rapproche à une valeur réelle par le rapport de deux entiers. Conduit à des opérations arithmétiques qui sont difficiles à implémenter.

Nous nous limitons aux représentations à virgule fixe et à virgule flottante.

#### I.3.1 Format à virgule fixe

Les données en format virgule fixe sont représentées comme étant des nombres fractionnaires à virgule fixe (exemple -1.0 à 1.0), ou comme des entiers classiques. Il existe deux types pour ce format, format virgule fixe pour les entiers non signés et le format virgule fixe pour les entiers signés.

##### I.3.1.1 Format à virgule fixe : entiers non signés

La représentation d'un entier  $x$  est notée par le vecteur suivant :

$$x = (x_{n-1}, x_{n-2}, \dots, x_1, x_0) \quad (I.1)$$

L'indexation utilisée à l'origine zéro est concerné le digit le plus à gauche.

Un système de numération pour représenter  $x$  se compose des éléments suivants :

- ✓ Le nombre de chiffres  $n$ .
- ✓ Un ensemble de valeurs numériques pour les chiffres. Nous appelons  $D_i$  l'ensemble des valeurs des  $x_i$ .

Le cardinal de  $D_i$  est dénoté par l'ensemble de valeurs  $|D_i|$  de  $x_i$ . Par exemple,  $\{0, 1, 2, \dots, 9\}$  est le chiffre fixé pour les systèmes de numération décimal classiques avec cardinalité 10.

- ✓ Une règle d'interprétation qui correspond à un mappage entre l'ensemble des vecteurs-chiffres et l'ensemble de valeurs des entiers.

L'ensemble des entiers, chacun représenté par un vecteur à n chiffres.

Pour cela la représentation de nombre est comme suit :

$$x = \sum_{i=0}^{n-1} x_i \times w_i \quad (\text{I.2})$$

où  $w = (w_{n-1}, w_{n-2}, \dots, w_1, w_0)$  est le vecteur-poids.

Dans laquelle le vecteur-poids est lié au vecteur- Radix  $r = (r_{n-1}, r_{n-2}, \dots, r_1, r_0)$  comme suit :

$$w_0 = 1; \quad w_i = w_{i-1} \times r_{i-1} \quad (1 \leq i \leq n-1) \quad (\text{I.3})$$

Cela équivaut à :

$$w_0 = 1; \quad w_i = \prod_{j=0}^{i-1} r_j \quad (\text{I.4})$$

Les Systèmes de numération de base sont classés selon le vecteur-radix en systèmes fixe-radix et en systèmes mixtes-radix.

Dans un système de radix-fixe, tous les éléments du vecteur-radix ont la même valeur r (la base). Par conséquent, le vecteur-poids comme suit :

$$w = (r^{n-1}, r^{n-2}, \dots, r^2, r, 1) \quad (\text{I.5})$$

Et les jeux de chiffres sont :

$$D_i = D \quad (1 \leq i \leq n-1) \quad (\text{I.6})$$

et

$$x = \sum_{i=0}^{n-1} x_i \times r^i \quad (\text{I.7})$$

Les bases les plus fréquemment utilisées sont des puissances de deux, par exemple 2 (binaire), 4 (quaternaire), 8 (octal), 16(hexadécimal), et ainsi de suite. La gamme correspondante de x représentée avec n chiffres radix-r est :

$$0 \leq x \leq r^n - 1 \quad (\text{I.8})$$

Selon l'ensemble des données numériques (Di), les systèmes de numération de base sont classés en des systèmes redondants et des systèmes non redondants.

Un système de numération est non redondant si chaque vecteur-chiffre représente un autre nombre entier ; autrement dit, si le mappage de représentation est un à un.

Il est redondant s'il existe des entiers qui sont représentée par plus d'un vecteur-chiffre.

### I.3.1.2 Format à virgule fixe : entiers signés

Pour les entiers signés (positifs et négatifs) deux représentations courantes :

- **La représentation du signe et magnitude (SM)** : en SM, un entier signé  $x$  est représenté par un couple  $(x_s, x_m)$ , où  $x_s$  est le signe et la  $x_m$  est la magnitude (entier positif).

Les deux valeurs de signe (+, -) sont représentées par une variable binaire, « 0 » correspond à un plus (+) et « 1 » pour un moins (-).

Si on utilise un système de radix- $r$  classique, la gamme des entiers, pour  $n$  chiffres signés, est  $0 \leq x_m \leq r^n - 1$ . Notez que zéro a deux représentations :  $x_s = 0, x_m = 0$  (zéro positif) et  $x_s = 1, x_m = 0$  (zéro négatif).

- **La représentation par un complément** : dans ce système, il n'y a pas de séparation entre la représentation du signe et la représentation de l'ampleur, mais l'ensemble entier signé est représenté par un entier positif. Les représentations des entiers positifs sont appelées formes véritables et ceux des entiers négatif, les formes de complément. Le complément est exprimé dans la cas général en radix- $r$ , nous ne considérons que le cas particulier de radix-2, appelé les deux représentations du complément.

Dans la représentation du complément à deux, un entier  $x$  signé (bit vecteur) avec  $n$  chiffres est représenté comme suit :

$$x = -x_{n-1} \times 2^{n-1} + \sum_{i=0}^{n-2} x_i \times 2^i \quad (\text{I.9})$$

Pour convertir un vecteur de bits d'une valeur, nous utilisons le fait que le bit le plus significatif ( $x_{n-1}$ ) de  $x$  a un poids négatif, tandis que les bits restants ont des poids positifs.

**Exemple :**

$$x = (11011) \rightarrow -16 + 8 + 2 + 1 = -5.$$

La représentation complément à deux (EQ I.9) possède les propriétés suivantes :

- La représentation sous forme de zéro est unique. Zéro est obtenue lorsque tous les chiffres sont mises à zéro.

- La gamme des nombres n'est pas symétrique puisque  $x = -2^{n-1}$  est représentable mais  $x = 2^{n-1}$  n'est pas représentable. La gamme est  $[-2^{n-1}, 2^{n-1} - 1]$ .

### I.3.1.3 Arithmétique virgule fixe pour le complément à deux

Avant de décrire les opérations arithmétiques, nous devons d'abord officialiser la représentation des nombres au format à virgule fixe [27].

- Soit  $\beta$  le nombre de chiffres (bits) d'un nombre signé (bit-vecteur)  $x$ , et  $\gamma$  le nombre de chiffres de la partie décimale de  $x$  (Fig. I.1).

- Soit  $\alpha$  le nombre de chiffres de la partie entière ( $\alpha = \beta - \gamma$ ). Enfin, FPR $x$  désigne le triple  $(\alpha, \beta, \gamma)$  définissant la représentation à virgule fixe de  $x$ . par conséquent,  $x = x_{Int} + x_{fr}$  s'écrit

$x = (x_{\alpha-1} \dots x_1 x_0 \bullet x_{-1} \dots x_{-\gamma+1} x_{-\gamma})$ , telle que

$$x = -x_{\alpha-1} \times 2^{\alpha-1} + \sum_{i=-\gamma}^{\alpha-2} x_i \times 2^i \quad (\text{I.10})$$

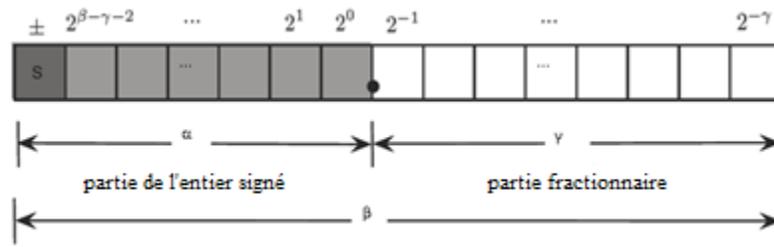


Figure I.1 format virgule fixe pour les nombres signés en complément à deux[8].

Pour convertir un nombre réel x FPR, nous procédons comme suit :

- B est une valeur donnée, puisqu'elle doit être égale à la largeur en bits du chemin de données.
- Le nombre de bits de la partie entière est donné par :

$$\alpha_x = \lfloor \log_2 |x| \rfloor + 2 \tag{I.11}$$

Avec  $\lfloor a \rfloor$  est la fonction floor qui arrondit au nombre entier inférieur ou égal le plus proche à x. par conséquent, nous déterminons la partie fraction comme  $\gamma_x = \beta_x - \alpha_x$ , et le FPR de x est donné par :  $FPR_x = (\beta_x, \alpha_x, \gamma_x)$

Par contre, dans la représentation en virgule fixe, x est représentée par l'entier  $N_x$  tels que :

$$N_x = \lfloor x \times 2^{\gamma_x} \rfloor \tag{I.12}$$

où  $\lfloor a \rfloor$  est la fonction qui arrondit au nombre entier le plus proche. Ainsi, le x est approximé par

$$x^\dagger = N_x \times 2^{-\gamma_x} \tag{I.13}$$

L'addition et la multiplication sont les deux opérations arithmétiques les plus utilisées dans les systèmes linéaires. En format virgule fixe, elles sont gérées comme suit.

### I.3.1.4 La Multiplication en virgule fixe

Considérons l'opération  $z = x \times y$ , avec  $(\beta_x, \alpha_x, \gamma_x)$  et  $(\beta_y, \alpha_y, \gamma_y)$ . FPRz est donnée par :

$$FPR_z = (\beta_x + \beta_y, \alpha_x + \alpha_y, \gamma_x + \gamma_y) \tag{I.14}$$

et l'opération de multiplication est réalisée par  $N_z \leftarrow N_x \times N_y$ . Dans l'équation 3.14 la multiplication est effectuée en double précision depuis  $\beta_z = \beta_x + \beta_y$ , mais généralement le bit-largeur ( $W_{dp}$ ) du chemin de données- est plus petit que  $(\beta_x + \beta_y)$ . Dans ce cas, laissez-nous désigner  $\beta_{op} = W_{dp}$  et donner l'expression générale de FPRz :

$$FPR_z = (\beta_{op}, \alpha_x + \alpha_y, \beta_{op} - \lfloor \alpha_x + \alpha_y \rfloor) \tag{I.15}$$

et l'opération est réalisée par

$$N_z \leftarrow (N_x \gg s_x) \times (N_y \gg s_y) \tag{I.16}$$

où  $s_x$  et  $s_y$  et le bit droit déplacements appliqués sur  $N_x$  et  $N_y$  telle que  $s_x + s_y = (\beta_x + \beta_y) - \beta_{op}$ .

Mais Si  $\beta_{op} = (\beta_x + \beta_y)$ ,  $s_x = s_y = 0$ . Les cas spéciaux et généraux, données par l'équation I.14

et I.15, respectivement, sont illustrés par la figure I.2. Un nombre de bits  $S_{cst}$  est tronqué de l'opérande Y de telle sorte que  $\beta_{op} = W_{dp}$ .

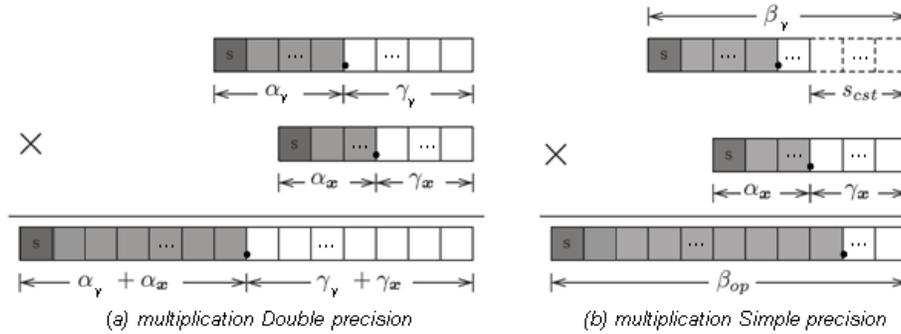


Figure I.2 (a) Double et (b) simple précision[8]

### I.3.1.5 Addition en virgule fixe :

L'addition de deux nombres à virgule fixe nécessite la même position du point. La précision complète est donnée par :

$$\begin{cases} \alpha_z = \max(\alpha_x + \alpha_y) + 1 \\ \gamma_z = \max(\gamma_x + \gamma_y) \\ \beta_z = \alpha_z + \gamma_z \end{cases} \quad (I.17)$$

Le cas général (précision limitée) où est donnée par :

$$\begin{cases} \alpha_z = \max(\alpha_x + \alpha_y) + 1 \\ \gamma_z = \beta_{op} - \max(\alpha_x + \alpha_y) - 1 \\ \beta_z = \beta_{op} \end{cases} \quad (I.18)$$

et l'opération d'addition est réalisée par

$$N_z \leftarrow (N_x \gg s_x) \times (N_y \gg s_y) \quad (I.19)$$

**Avec**  $s_x = \gamma_x - \gamma_z$  et  $s_y = \gamma_y - \gamma_z$

### I.3.1.6 Détection d'un Débordement (overflow)

Dans la représentation du complément à deux, un dépassement de capacité se produit lorsque les opérandes sont de même signe et le résultat de l'addition représente un nombre entier de signe différent. Étant donné que le signe est déterminé par le bit le plus significatif ( $x_{n-1}$ ), la détection de dépassement de capacité est spécifiée par l'expression suivante :

$$AVF = (\bar{x}_{n-1} \times \bar{y}_{n-1} \times z_{n-1}) + (x_{n-1} \times y_{n-1} \times \bar{z}_{n-1}) \quad (I.20)$$

### I.3.2 Format à virgule flottante

La représentation en virgule flottante est utilisée pour la représentation des nombres réels. Comme dans la représentation en virgule fixe, la représentation en virgule flottante est codée sur un nombre fini de bits. Pour un système spécifique à virgule flottante, un nombre réel qui est exactement représenté dans le système est appelé un nombre en virgule flottante. Dans le cas d'un dépassement de capacité positifs et négatifs les nombres sont représentés par des nombres à virgule flottante qui ont une valeur qui se rapproche du nombre réel. Le processus de rapprochement est appelé arrondissement et génère une erreur.

Un nombre  $x$  en virgule flottante est représenté par un triple  $(S_x, M_x, E_x)$ , tels que :

$$x = (-1)^{S_x} \times M_x \times b^{E_x} \quad (I.21)$$

où  $b$  est une constante appelée la base ;  $E_x$  est un exposant entier signé ;  $M_x$  est la mantisse ; et  $S_x \in \{0,1\}$  est le signe de la mantisse.

#### I.3.2.1 La gamme dynamique

Le but d'utiliser la représentation à virgule flottante est d'augmenter la gamme dynamique, par rapport à une représentation en virgule fixe. Cette gamme dynamique est définie comme le rapport entre le nombre le plus grand et le nombre le plus petit non nuls et positif qui peut être représentée [1]. Pour obtenir une représentation en virgule fixe à l'aide de digits de  $n$  radix- $r$ , la gamme dynamique est :

$$DR_{fxt} = r^n - 1 \quad (I.22)$$

En revanche, pour la représentation à virgule flottante :

$$DR_{fpt} = \frac{M_{\max} \times b^{E_{\max}}}{M_{\min} \times b^{E_{\min}}} \quad (I.23)$$

Par exemple, si les chiffres de  $n$  sont partitionnés ainsi que  $m$  chiffres sont utilisés pour les chiffres mantisse et  $n-m$  pour l'exposant et  $b = r$ , nous obtenons :

$$DR_{fpt} = (r^m - 1) \times r^{(r^{n-m} - 1)} \quad (I.24)$$

Par exemple, si  $n = 32$ ,  $m = 24$  et  $r = 2$ , les plages dynamiques correspondantes sont :

$$DR_{fxpt} = 2^{32} - 1 \approx 4.3 \times 10^9$$

$$DR_{fpt} = (2^{24} - 1) \times 2^{2^8 - 1} \approx 9.7 \times 10^{83}$$

Une large gamme dynamique est requise dans de nombreuses applications pour éviter les dépassements de capacité positifs et négatifs. Si la gamme dynamique de la représentation à virgule fixe  $n$  n'est pas suffisante, des opérations complexes doivent être incluses dans le programme. Ainsi, un système à virgule flottante est préférable dans de telles applications.

### I.3.2.2 Précision

Dans la représentation en virgule flottante, Pour la même plage de valeurs représentables, virgule flottante a tendance à être mieux que virgule fixe en termes de l'erreur moyenne.

Bien que la représentation en virgule flottante fournisse une grande gamme dynamique et une précision supérieure à la représentation en virgule fixe, elle est beaucoup plus coûteuse à mettre en œuvre. Dans les applications des systèmes embarqués tolérant un certain degré d'imprécision, la notation à virgule fixe est la plus utilisées pour augmenter le débit et diminuer la superficie, le retard, et l'énergie.

## I.4 Les systèmes de représentation des nombres

Les systèmes de numération sont développés afin de permettre une réduction pour la complexité des opérations arithmétiques. Chaque système de numération présente des propriétés numériques spécifiques, les opérations arithmétiques sont gérées différemment parmi ces systèmes de numération on peut trouver.

### I.4.1 La représentation numérique canonique des nombres signés (CSD)

Est la forme de la représentation canonique des entiers signés développée par Avizienis [11,8] en 1961.

Est un radix-fixe redondant ( $r = 2$ ) sa représentation définie comme suit :

$$x = \sum_{i=0}^{n-1} x_i \times 2^i \quad \text{avec } x_i \in D = \{\bar{1}, 0, 1\} \text{ et } \bar{1} = -1 \quad (I.25)$$

L'avantage principal de ce système est que les opérations d'addition/soustraction peuvent être faites sans rapport de propagation, et avec un traitement accéléré, notamment pour les grands opérandes.

Dans le cas des entiers signés, l'entier « sept », par exemple, possède plusieurs représentations :  
 $(0111)_2 = 4 + 2 + 1 = 7$ ,  $(10\bar{1}1)_2 = 8 - 2 + 1 = 7$ ,  $(1\bar{1}11)_2 = 8 - 4 + 2 + 1 = 7$ ,  $(100\bar{1})_2 = 8 - 1 = 7$

Le nombre rationnel « 5/8 », par exemple ; peut-être écrit différemment :

$$\frac{5}{8} = 0.625 = (0.101)_2 = (1.\bar{1}01)_2 = (1.\bar{1}1\bar{1})_2 = (1.0\bar{1}\bar{1})_2$$

La conversion de la SD en CSD se fait comme suit :

- ✓ Convertir une longue chaîne de 1.
- ✓ Fusionner des chiffres adjacents de signes opposés : en utilisant le format de la CSD sur une valeur de  $n$  bits, le nombre de chiffres différent de zéro est délimité par les  $(n + 1) / 2$  et il tend asymptotiquement à une valeur moyenne de  $(n/3) + (1/9)$ . Par rapport à la représentation binaire traditionnelle nécessitant  $n/2$  chiffres sur la moyenne, CSD permet une économie de 33 % en chiffre différent de zéro. Cela signifie que, dans la multiplication constante, 33 % moins d'additions/soustractions sont nécessaires, ce qui conduit à beaucoup compact implémentations de systèmes LTI [2]. C'est la raison pour laquelle les CSD est populaire.

Malgré le fait que le CSD minimise le nombre de chiffres différent de zéro pour la représentation de la constante, c'est loin d'être optimale. Il est possible de décomposer la valeur de  $n$  bits pour réduire davantage le nombre d'opérations. Ceci peut être réalisé à l'aide de plusieurs systèmes de nombre complexe.

#### I.4.2 Système de numération à double Base (DBNS)

L'arithmétique DBNS a été développé par Dimitrov en 1999 [12]. Dans DBNS, un entier  $x$  est exprimé à l'aide des bases 2 et 3, comme suit :

$$x = \sum_{i,j} d_{i,j} \times 2^i \times 3^j \quad \text{Avec } d_{i,j} \in D = \{0,1\} \quad (\text{I.26})$$

Par exemple, le nombre entier  $x = (10603)_{10}$  est écrit en DBNS comme suit :

$$x = (3^2 \times 2^8) + (3 \times 2^5) + (3^0 \times 2^{13}) + (3^0 \times 2^3) + (3^1 \times 2^0) = (10603)_{10}.$$

Selon EQ. I.2 et I.3, La représentation en DBNS est extrêmement redondante. La représentation canonique de DBN (CDBNR) qui exprime un entier donné comme une somme d'un nombre minimal de 2-nombres entiers est très complexe à déterminer (problème NP-complet). Ainsi, les opérations arithmétiques dans ce système de numération ne garantissent pas que les résultats qui sont obtenus soient sous la forme minimale.

#### I.4.3 Système de numération de résidus (RNS)

Le concept de RNS remonte à 1500 ans en Chine [4]. Un nombre  $x$  est représenté par le vecteur de ses résidus en ce qui concerne les premiers modules  $k \ m_{k-1} > m_{k-2} > \dots > m_1 > m_0$ . Le résidu  $x_i$  de  $x$  à l'égard de l' $i$ -ème module  $m_i$  est semblable à un digit, et la représentation de l'ensemble  $k$ -résidu de  $x$  peut être considérée comme un  $k$ -digits. Sans, nous écrire

$$x = (x_{k-1} | x_{k-2} | \dots | x_1 | x_0)_{RNS(m_{k-1}|m_{k-2}|\dots|m_1|m_0)}$$

$$x_i = x \bmod m_i = \langle x \rangle_{m_i} \quad \text{avec } x_i \in \{0,1,2,\dots, m-1\} \quad (\text{I.27})$$

Le vecteur des modulus  $RNS(m_{k-1}|m_{k-2}|\dots|m_1|m_0)$  peuvent être supprimés de l'indice lorsque nous nous sommes mis d'accord sur un ensemble par défaut. Le produit  $M$  des premiers modulus

de  $k$  est le nombre de différentes valeurs représentables dans le RNS et est connu comme sa gamme dynamique :

$$M = m_{k-1} \times m_{k-2} \times \dots \times m_1 \times m_0 \quad (I.28)$$

Par exemple,  $M = 8 \times 7 \times 5 \times 3 = 840$  est le nombre total de valeurs distinctes qui sont représentables dans  $RNS = (8 | 7 | 5 | 3)$ . À cause de l'égalité

$$\langle -x \rangle_{m_i} = \langle M - x \rangle_{m_i} \quad (I.29)$$

Les 840 valeurs disponibles peuvent être utilisées pour représenter les nombres de 0 jusqu'à 839, ou de -420 jusqu'à 419, ou tout autre intervalle de 840 entiers consécutifs. En effet, les nombres négatifs sont représentés à l'aide d'un système de complément avec la constante de complémentation  $M$ . Voici quelques exemples de chiffres dans  $(8 | 7 | 5 | 3)$  :

$(0 | 0 | 0 | 0)_{RNS}$  Représente 0 or 840 or ...

$(1 | 1 | 1 | 1)_{RNS}$  Représente 1 or 841 or ...

$(2 | 2 | 2 | 2)_{RNS}$  Représente 2 or 842 or ...

$(0 | 1 | 3 | 2)_{RNS}$  Représente 8 or 848 or ...

$(5 | 0 | 1 | 0)_{RNS}$  Représente 21 or 861 or ...

$(0 | 1 | 4 | 1)_{RNS}$  Représente 64 or 904 or ...

$(2 | 0 | 0 | 2)_{RNS}$  Représente -70 or 770 or ...

$(7 | 6 | 4 | 2)_{RNS}$  Représente -1 or 839 or ...

Représentation RNS n'est pas redondante au sein de l'intervalle choisi par 840 entiers consécutifs.

### I.4.3.1 Addition et Multiplication

Le signe d'un nombre RNS peut être changé en complétant indépendamment chacun de ses chiffres en ce qui concerne son modulo. De même, l'addition, la soustraction, et la multiplication peuvent être effectuées en agissant indépendamment sur chaque chiffre. Les exemples suivants pour RNS  $(8 | 7 | 5 | 3)$  illustrent le processus d'addition, soustraction et multiplication, respectivement :

$(5 | 5 | 0 | 2)$  RNS représente  $x = +5$

$(7 | 6 | 4 | 2)$  RNS représente  $y = -1$

$(4 | 4 | 4 | 1)$  RNS  $x + y$  ;, etc.

$(6 | 6 | 1 | 0)$  RNS  $x - y$  ;, etc.

$(3 | 2 | 0 | 1)$  IA  $x \times y$  ;, etc.

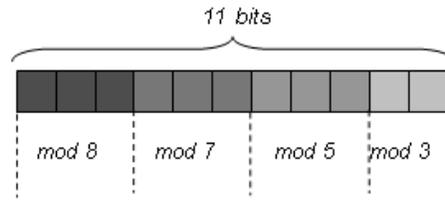


Figure I.3 Format RNS pour le nombre (8/7/5/3)[8]

La vitesse et la simplicité sont les principaux avantages des opérations arithmétiques. Dans le cas de l'addition, par exemple, la propagation de retenue est limitée à un seul résidu (quelques segments). Ainsi, la représentation RNS a résolu à peu près le problème de propagation de la retenue. En ce qui concerne la multiplication, un multiplicateur de 4 x 4 par exemple, est quatre fois plus simple qu'un multiplicateur de 16 x 16, en plus d'être beaucoup plus rapide. En fait, étant donné que les résidus sont légers (par exemple 6 bits de largeur), il est tout à fait possible d'appliquer l'addition, la soustraction, et la multiplication par recherche dans la table directement. Avec des résidus de 6 bits, chaque opération nécessite une table de 4 Kx6. Ainsi, à l'exception de la division, un module d'unité arithmétique complète pour un résidu de 6 bits peut être implémenté avec 9 Ko de mémoire.

Malheureusement, ce que nous gagnons en ce qui concerne la vitesse et la simplicité d'implémentation de l'addition, la soustraction, et la multiplication, peut être réduit à néant par la complexité de la division et la difficulté de certaines opérations auxiliaires telles que le test du signe, Comparaison de magnitude et détection de dépassement de capacité.

#### I.4.4 Système de numération radix-2<sup>r</sup>

La représentation de radix-2<sup>r</sup> a été développée par Sam en 1990 [13]. en radix-2<sup>r</sup>, le complément à deux d'un nombre x est comme suit :

$$\begin{aligned}
 x &= \sum_{j=0}^{(n/r)-1} (x_{rj-1} + 2^0 x_{rj} + 2^1 x_{rj+1} + 2^2 x_{rj+2} + \dots + 2^{r-2} x_{rj+r-2} - 2^{r-1} x_{rj+r-1}) \times 2^{rj} \\
 &= \sum_{j=0}^{(n/r)-1} Q_j \times 2^{rj} \quad \text{Avec } Q_j \in D = \{-2^{r-1}, -2^{r-1} + 1, \dots, -1, 0, 1, \dots, 2^{r-1} - 1, -2^{r-1}\} \quad (I.30)
 \end{aligned}$$

Quand  $x_{-1} = 0$  et  $r \in \mathbb{N}^*$ . Pour une simplicité et sans perte de généralité, nous supposons que r est un diviseur de N. En EQ. I.30, de complément à deux de x est divisée en n/r complément à deux segments, chacune des r + 1 bit de longueur. Chaque paire de deux segments continues a un bit qui se chevauchent.

En fait, la représentation SD (EQ. I.25) est un cas particulier de représentation radix-2r (EQ. I.30) pour r = 1.

Le signe de l'expression est donné par le bit  $x_{rj+r-1}$  et, avec  $k_j \in \{0, 1, 2, \dots, r-1\}$  et  $m_j \in OM(2^r) \cup \{0\}$  où  $OM(2^r) = \{1, 3, 5, \dots, 2^{r-1} - 1\}$ .  $OM(2^r)$  est l'ensemble des chiffres positifs impairs radix-2r recodage, avec  $|OM(2^r)| = 2^{r-2}$ .  $|Q_j| = 0$  Correspond  $m_j = 0$ . Enfin, x peut être exprimée comme suit :

$$x = \sum_{j=0}^{(n/r)-1} (-1)^{x_{rj+r-1}} \times m_j \times 2^{rj+k_j} \tag{I.31}$$

**Exemple illustratif :**

Le produit 10599 x X est exprimé en CSD, DBNS, et RADIX-2<sup>r</sup>. Pour représenter ce produit en RADIX-2<sup>r</sup>, il est d'abord nécessaire de représenter (10599)<sub>10</sub> en complément à deux, ce qui donne (010100101100111)<sub>2</sub>. Donc dans la représentation complément à deux la longueur de bits de la constante sera N+1 (14+1=15 pour 10599).

Pour N = 14 correspond à r = 3 pour C= 10599 les équations (I.30) et (III.3) devient respectivement :

$$C = \sum_{j=0}^4 Q_j \times 2^{3j} \text{ Et } P_{RADIX} = \sum_{j=0}^4 (-1)^{c_{3j+2}} \times (m_j \times X) \times 2^{3j+k_j}$$

$$C = \sum_{j=0}^4 (c_{3j-1} + 2^0 c_{3j} + 2^1 c_{3j+1} + 2^2 c_{3j+2}) \times 2^{rj}$$

La table de vérité équivalant à cette représentation est calculée en utilisant l'équation précédent, les résultats de calcul des m<sub>j</sub>, k<sub>j</sub> illustré dans le tableau I.1.

Tableau I.1 La table de verité pour 10599 en RADIX-2r

Q <sub>j</sub>				m <sub>j</sub>	k <sub>j</sub>
C <sub>3j+2</sub>	C <sub>3j+1</sub>	C <sub>3j</sub>	C <sub>3j-1</sub>		
0	0	0	0	0	0
0	0	0	1	1	0
0	0	1	0	1	0
0	0	1	1	1	1
0	1	0	0	1	1
0	1	0	1	3	0
0	1	1	0	3	0
0	1	1	1	1	2
1	0	0	0	1	2
1	0	0	1	3	0
1	0	1	0	3	0
1	0	1	1	1	1
1	1	0	0	1	1
1	1	0	1	1	0
1	1	1	0	1	0
1	1	1	1	0	0



Tableau I.2 nombre des chiffres requis pour chaque système pour la valeur  $(10599)_{10}$ 

Systeme de numérisation	Expression arithmétique	Nombre des chiffres
Binary	$x = 2^{13} + 2^{11} + 2^8 + 2^6 + 2^5 + 2^2 + 2^1 + 2^0$	8
SD	$x = 2^{13} + 2^{12} - 2^{11} + 2^9 - 2^8 + 2^7 - 2^5 + 2^3 - 2^0$	9
CSD	$x = 2^{13} + 2^{11} + 2^9 - 2^7 - 2^5 + 2^3 - 2^0$	7
DBNS	$x = (3^2 \times 2^8) + (3 \times 2^5) + (3^0 \times 2^{13}) + (3^0 \times 2^3) - (3^0 \times 2^0)$	$5+2=7$
RNS	$x$ is represented by $(7 \mid 28 \mid 9)_{\text{RNS}(32 \mid 31 \mid 15)}$	$3+3+2=8$
Radix- $2^r$	$x = (3 \times 2^{12}) - (7 \times 2^8) + (3 \times 2^5) + (7 \times 2^0)$	$4+2=6$

Note que les termes  $3$ ,  $3^2$ , et  $7$  dans DBNS et RADIX- $2^r$  considéré comme des extra-chiffres

Comme nous l'avons déjà mentionné au début de ce chapitre, l'objectif d'un système de numération est l'absorption maximale de logique-redondante dans la représentation des nombres. Cela conduira à moins de digits et par conséquent à moins de ressources matérielles et plus de vitesse. Ce processus d'absorption peut être simple, par exemple dans le cas de CSD, ou très compliqués tels que dans le cas de systèmes DBNS et RNS. Selon le type de fonction arithmétique à mettre en œuvre, le processus d'absorption s'effectue soit en logiciel ou en matériel (hardware). Par exemple, dans la multiplication par une constante ( $C \times X$ ), le processus de réduction est réalisé en logiciel, alors que dans la multiplication par une variable ( $X \times Y$ ), la réduction est implémentée en matériel. Par conséquent, l'avantage de réduire le nombre de digits peut être compensé par l'effort d'absorption, nécessitant trop de temps-de calcul (logiciel), ou trop de ressources matérielles (tableau I.3).

Tableau I.3 Principales caractéristiques des systèmes numériques

Number System	Weighted System	Radix System	Fixed/Mixed Radix System	Canonical Form	Ease of use & Implementation	Frequently Used	Hardware Optimization
Binary	Yes	Yes	Fixed	Yes	A	A	E
SD	Yes	Yes	Fixed	Yes	B	A	D
CSD	Yes	Yes	Fixed	Yes	C	A	C
DBNS	Yes	No	–	U*	E	C	B
RNS	Yes	Yes	Mixed	U*	F	D	D
Radix- $2^r$	Yes	Yes	Fixed	U*	D	B	A

A-F : A, le meilleur ; F, le plus mauvais ; U : Inconnu ; \* : l'existence de la forme canonique doit être prouvée.

**I.5 Conclusion**

L'arithmétique binaire est un vaste sujet. C'est pourquoi nous avons volontairement limité notre étude selon notre besoin pour les principales idées et les concepts essentiels impliqués dans la conception des systèmes linéaires. Nous nous sommes concentrés surtout sur l'arithmétique à virgule fixe et les principaux systèmes de nombre qui lui sont liés. Une attention particulière a été consacrée à l'addition et la multiplication qui sont les deux principales opérations de toute architecture de système linéaire.

# **Chapitre II :**

## **Multiplication par une constante SCM/MCM**

## II.1 Introduction

Ce chapitre aborde le problème d'optimisation des ressources nécessaire pour les systèmes linéaires invariants (LTI), en se concentrant principalement au niveau de la multiplication par une constante simple/multiple (SCM/MCM) afin d'élaborer des solutions communes.

## II.2 Les systèmes LTI

Un système linéaire est un modèle mathématique basé sur des opérations linéaires. Une opération linéaire se conforme par deux propriétés, à savoir l'additivité et l'homogénéité. Étant donné deux vecteurs  $x$  et  $y$  et un scalaire  $c$ , ces Propriétés sont décrites comme suit :

$$\checkmark \text{ Additivité : } f(x + y) = f(x) + f(y) \quad (\text{II.1})$$

$$\checkmark \text{ Homogénéité : } f(c \times x) = c \times f(x) \quad (\text{II.2})$$

### II.2.1 Le rôle des systèmes linéaires invariants

Les Systèmes linéaires ont des rôles comme des abstractions mathématiques pour des modèles de calcul dans un grand nombre d'applications, dont principalement : la théorie du contrôle automatique, traitement du signal et télécommunications.

### II.2.2 Formulation des systèmes LTI :

Si  $X$  et  $Y$  sont des entrées et des sorties d'un type vecteurs, respectivement, et  $C$  est une matrice de transformation, le système LTI peut être décrit comme suit :

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} & \dots & C_{1n} \\ C_{21} & C_{22} & \dots & C_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ C_{m1} & C_{m2} & \dots & C_{mn} \end{bmatrix} \times \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \quad (\text{II.3})$$

Telque la matrice de transformation  $C$  est une matrice  $m \times n$ , où le  $C_{ij}$  représente la constante, le  $Y_i$  est le signal de sortie de produit de la ligne  $i$  de la matrice de transformation  $C$  et les échantillons d'entrée  $n$  de  $x$  :

$$Y_j = \sum_{j=1}^n C_{ij} \times X_j \quad (\text{II.4})$$

### II.3 La contrainte de la multiplication par une constante

Le problème de la mise en œuvre de la multiplication par une constante est apparu dans les années 1970. Plusieurs microprocesseurs dans ce temps (tels que les Intel 8008) n'a pas eu des multiplieurs, d'où la multiplication doit être fait avec des additions, soustractions et des décalages. Et même lorsque l'instruction de la multiplication est devenue disponible, l'exécution nécessite généralement plusieurs cycles d'horloge.

Aujourd'hui, et avec la notion de pipeline la résolution de problème de la multiplication par une constante pourrait conduire à une réduction pour le temps d'exécution mais pas comme une solution satisfiable.

SCM résout uniquement le problème de la multiplication d'une variable ( $X_j$ ) par une constante ( $C_{ij}$ ). Pour être efficacement gérées, la mise en œuvre matérielle du  $C_{ij} \times X_j$  doit être faite exclusivement par des additions, des soustractions et des décalages.

Afin de conduire à des avantages considérables en ce qui concerne la durée d'exécution, le domaine, débit et puissance, la complexité du calcul de SCM semble encore inconnue. C'est seulement conjecture (non preuve) d'être NP-difficile [14,8]. Mais parce que la solution-espace à explorer est énorme, il faut utiliser l'heuristique de minimisation de nombres des additions.

#### Exemples illustratifs :

L'exemple est  $45 \times X_j$ , de nombreuses solutions possibles, seulement quatre des solutions sont présentées :

$$45 \times X_j = (101101)_2 \times X_j = X_j \times 2^5 + X_j \times 2^3 + X_j \times 2^2 + X_j \quad (\text{II.5})$$

$$45 \times X_j = (63-18) \times X_j = [(111111)_2 - (010010)_2] \times X_j = (X_j \times 2^6 - X_j) - X_j \times 2^4 - X_j \times 2 \quad (\text{II.6})$$

$$45 \times X_j = (10\bar{1}0\bar{1}01)_2 \times X_j = X_j \times 2^6 - X_j \times 2^4 - X_j \times 2^2 + X_j \quad (\text{II.7})$$

$$45 \times X_j = U \times 2^2 + U \text{ avec } U = X_j \times 2^3 + X_j \quad (\text{II.8})$$

- EQ. II.5 est la méthode la plus simple pour la multiplication à simple constante. elle transforme le constant  $C_{ij}$  dans sa représentation binaire, convertit les 1 en décalage basés sur leurs positions, et additionne les valeurs décalées.  
Pour  $C_{ij} = 45$ , EQ. II.5 nécessite trois additions et trois opérations de décalage. Le nombre des additions à l'aide de la représentation binaire est inférieur au nombre d'instances « 1 ».
- autre méthode (EQ. II.6) utilise les décalages et les soustractions en traduisant les valeurs « 0 » en décalage tout en les soustrayant de la constante de la même longueur, composée seulement des 1. La constante de 45 requiert six bits et la constante correspondante de six bits de tous les 1 (111111) est 63. Le terme  $(X_j \times 2^6 - X_j)$  représente le nombre de 63 et suit avec deux méthodes,  $X_j \times 2^4$  et  $X_j \times 2$  représentent les termes 16 et 2, respectivement ( $16 + 2 = 18$ ).
- La représentation de la CSD (EQ. II.7) encode un nombre constant en utilisant le nombre minimal de chiffres différent de zéro. Par conséquent, lors de la transformation d'une multiplication constante en une séquence de décalage et additions, la représentation de la

CSD donne le nombre minimal des additions. Mais dans ce cas spécial ( $C_{ij} = 45$ ), ne fournit aucun avantage sur EQ. II.5 et II.6.

- EQ. II.8 permet d'obtenir le nombre minimal d'addition pour  $C_{ij} = 45$ . La réduction des additions provient de la mise en commun du terme  $U = 9j \times X$ . Cela est bien illustré par la figure. II. 1 b, où les nœuds gris signalent une opération de décalage, et les rouges indiquent une addition. Le nombre total d'opérations est deux additions. Notez que plusieurs solutions avec 2 additions peuvent exister (Fig. II. 1 a, II. 1 b et c de II.1).

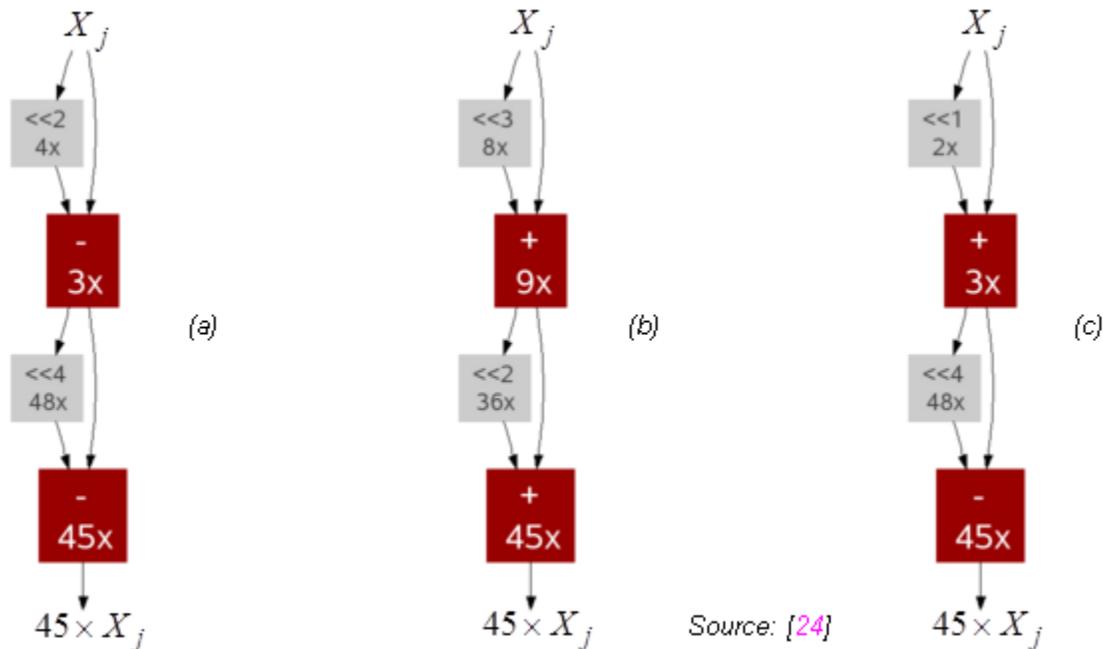


Figure II.1 le nombre minimal d'additions pour  $45j \times X$ .

Les solutions sont données par le site web de spirale ([www.spiral.net](http://www.spiral.net)). Le signe «  $\ll a$  » signifie un décalage d'une position ( $\times 2^a$ ).

L'heuristique SCM vise à fournir des solutions optimales dans un temps de calcul raisonnable. La prévisibilité est une autre caractéristique importante de l'heuristique de la SCM. Selon la taille des bits de la constante, il permet de connaître à l'avance (avant la mise en œuvre), le nombre maximal d'additions et le couût formant le chemin critique (vitesse).

### II.3.1 Multiplication de multiple-Constant (MCM)

Envisager III.3 l'équation linéaire ; MCM est une extension du SCM, où est la seule variable  $X_j$  multiplié par l'ensemble des constantes  $C_{1j}, C_{2j}, C_{3j}, C_{mj}$  (la colonne entière  $j$ ). Définis comme tels, le Problème MCM est au moins aussi dur que le problème de la SCM, puisque ce dernier est un sous ensemble du problème SCM ; étant donné que le problème de la multiplication de single-constant a été conjecturé NP-difficile, le problème de la MCM est

également conjecturé NP-difficile. Une solution potentielle au problème MCM est simplement optimisée par SCM indépendamment.

### Exemple illustratif

Nous effectuons les deux multiplications suivantes :  $81 \times X_j$  et  $23 \times X_j$ . Fig. II. 2 a et II. 2 b représentent la meilleure optimisation pour chaque cas individuellement. Les nœuds indiquent plus, alors que le spectacle des bords du montant du transfert. Fig. II. 2 a et II. 2 b nécessitent deux additions, entraînant quatre additions à effectuer deux multiplications. Fig. II.2 c montre l'optimisation simultanée des deux variables. Les variables peuvent partager une multiplication commune  $9x$ , par conséquent, le nombre total des additions pour les deux variables est réduit d'une unité (soit un total de trois additions).

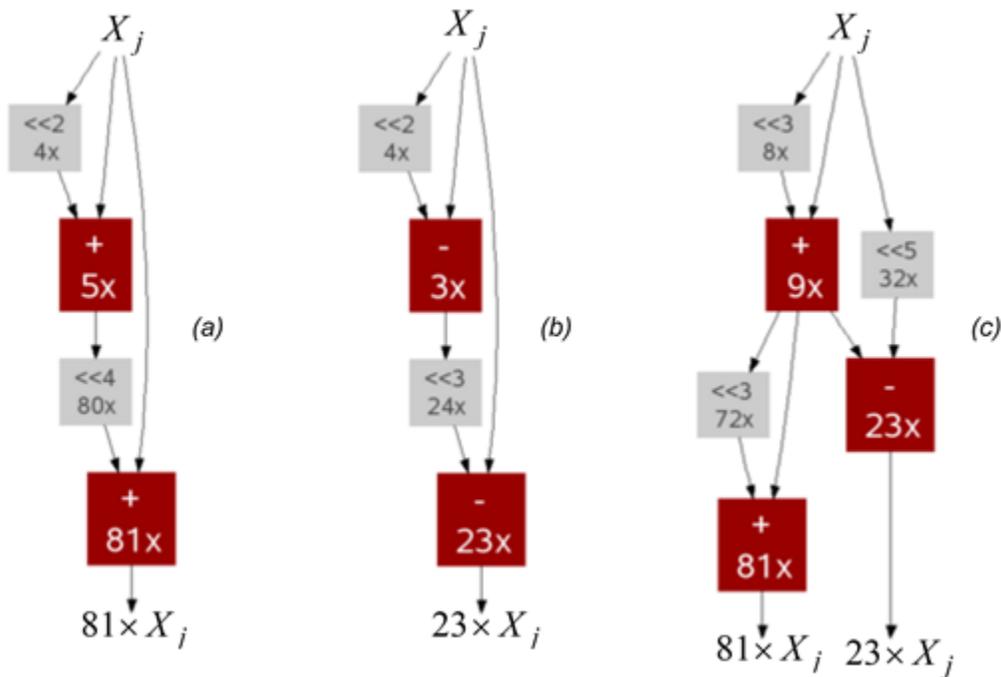


Figure II.2 Multiplication des constantes 81 et 23.

(a), (b) l'optimisation de chaque variable indépendamment ; Cela nécessite deux additions par constante pour un total de quatre additions. (c) l'optimisation simultanée de ces deux variables. Les variables peuvent partager une addition entraînant trois additions pour la multiplication de ces deux constantes.

### II.3.2 La formalisation du problème SCM

Avant de formaliser le problème de la SCM, nous devons définir clairement :

- Le type constant : positif, négatif, impair ou même. Presque toutes les propositions des heuristiques de SCM gérer seules les constantes impair/positif/négatif peuvent être dérivé simplement à l'aide des négation et des décalages. On note que dans ce cas un traitement extra des constantes est nécessaire.
- Les opérations autorisées : addition, soustraction, décalage vers la gauche, décalage vers la droite, et opération « ou ».

Limiter le nombre d'opérations rend le problème plus difficile à résoudre. La Plupart des projets des heuristiques permettent seulement l'addition, soustraction et décalage vers la gauche. Décalage à droite permet à quelques heuristiques [14,8] puisque il a des inconvénients, tout en apportant peu de valeur : par exemple, ils ne travaille pas avec l'arithmétique modulo  $2^k$ .

La définition formelle suivante s'applique à tous les types de constantes, c.-à-d., positif/négatif et impair/pair, permettant à des additions, soustractions et décalage à gauche seulement.

soit notre constant  $c$ ,  $c \in \mathbb{Z}$ . Une séquence finie des entiers signés  $u_0, u_1, u_2, \dots, u_q$  est acceptable pour  $c$  si elle satisfait les propriétés suivantes :

- Valeur initiale :  $u_0 = 1$  ;
- Pour  $i > 0$ ,  $u_i = s_i \times u_j \times 2^{a_i} + r_i \times u_k \times 2^{b_i}$  avec  $j, k < i$  ;  $s_i, r_i \in \{-1, 0, 1\}$  et  $a_i, b_i \in \mathbb{N}$  ;
- Valeur finale :  $u_q \times 2^{c_q} = c$  avec  $c_q \in \mathbb{N}$

Le problème est de trouver un algorithme qui prend le nombre  $C$  et qui génère une séquence acceptable  $(u_i)_{0 \leq i \leq q}$  aussi court que possible.  $q$  est appelé la qualité, ou la longueur. On note que pour retarder les déplacements, nous pouvons limiter  $b_i$  à 0 (cela casse la symétrie, mais fait moins variable). Ainsi, un nombre arbitraire  $X$  étant donné la solution correspondante itérativement calcule  $u_i \times X$  déjà calculé les valeurs  $u_j \times X$  et  $u_k \times X$ , jusqu'à obtenir  $c \times X$ . Note qu'avec cette formulation, on a calculé la valeur permise de réutiliser le tout déjà. C'est pourquoi certaines solutions générées peuvent besoin de stocker des résultats temporaires.

SCM/MCM est un problème fondamental dans le contrôle, DSP et les télécommunications

### II.3.3 Les algorithmes SCM/MCM existants

Les algorithmes MCM existants peuvent être divisés en quatre grandes classes :

- ✓ algorithmes de recodage des Chiffres tels que la représentation canonique chiffres signés (CSD) [11,8], stand recodage [15,8], et Dimitrov DBNS recodage [16,8] ;
- ✓ algorithme d'élimination des sous-expressions communes (CSE) à l'aide des critères spéciaux effectués après un premier recodage. Des exemples typiques sont Hartley [17,8], [26,8] de Lefèvre et Boullis [17,8] ;
- ✓ algorithme de réalisation des graphes acycliques (DAG) cette catégorie comprend les Bernstein [17,8], MAG [20,8] et Hcub [21,8] ;

- ✓ Algorithmes hybrides combinant CSE et DAG tel que l'algorithme optimal récent BIGE [14,8].

### II.3.4 Définition des mesures des algorithmes SCM/MCM

Les algorithmes SCM/MCM vise à produire une implémentation matérielle dans un laps de temps raisonnable en termes de la surface , vitesse et la consommation énergétique.

Étant donné une constante de N bits, les définitions principales des mesures susmentionnées sont :

**Définition 2.1** – la limite supérieure ( $U_{pb}$ ) : pour chaque constant  $C_i$  à N-bit correspond à  $A_i$  additions, pour l'implémentation de  $C_i \times X$  .  $U_{pb} = \max(A_i)$  .

**Définition de 2.2** – profondeur d'adder ( $A_{th}$ ) : soit  $D_i$  le nombre d'adder qui traversent l'entrée à une des sorties du circuit logique de la multiplication par une constante (le chemin le plus long).  $A_{th} = \max(D_i)$  .

**Définition 2.3** – moyenne ( $A_{vg}$ ) : pour chaque constante  $C_i$  à N-bit correspond à  $A_i$  additions, pour la mise en œuvre de  $C_i \times X$  ,  $A_{vg} = \sum_{i=1}^m A_i / m$  , où m est le nombre total des constantes.

En fait, bien que formellement, la profondeur de l'additionneur est définie par une estimation de chemin le plus long à travers le circuit logique ( chemin critique)[9,8].

### II.3.5 Les principales Limitations des algorithmes SCM/MCM existants :

L'optimisation des systèmes linéaires est un sujet essentiel qui a été au centre des efforts continus au cours des dernières années, ce qui entraîne un nombre impressionnant d'algorithmes SCM/MCM. Parmi les principales limitations de ces algorithmes existants on trouve :

- ✓ La Prévisibilité
- ✓ La durée de stockage en mémoire
- ✓ Le risque de débordement
- ✓ La facilité d'utilisation et la mise en œuvre

## II.4 Le Nouvel algorithme de recodage (RADIX-2<sup>r</sup>)

Il repose sur l'arithmétique radix-2<sup>r</sup> ,traite l'heuristique SCM/MCM susmentionné , facile à utiliser et plus sécurisé en cas d'un dépassement de capacité.

**Principe :**

Une constante C à N bits est exprimée en radix-2<sup>r</sup> comme suit :

$$C = \sum_{j=0}^{(n/r)-1} (C_{rj-1} + 2^0 C_{rj} + 2^1 C_{rj+1} + 2^2 C_{rj+2} + \dots + 2^{r-2} C_{rj+r-2} - 2^{r-1} C_{rj+r-1})$$

$$= \sum_{j=0}^{(n/r)-1} Q_j \times 2^{rj} \quad (\text{II.9})$$

Avec  $C_{-1}=0$  et  $r \in \mathbb{N}^*$

En cas d'un nombre signé; il s'agit de la formule suivante :

$$-2^{r-1}c_{rj+r-1} \times 2^{rj} + c_{rj+r-1} \times 2^{r(j+1)} = c_{rj+r-1} \times 2^{rj+r-1}.$$

La fonctionnalité d'algorithme RADIX- $2^r$  a une prévisibilité permet la génération des contrôleurs LTI entièrement capables de satisfaire les différentes exigences, telles que :

- Générer un contrôleur comprenant un nombre minimal d'additions (contrôleur plus compact);
- Générer un contrôleur avec un chemin critique plus court ( contrôleur plus rapide) ;
- Activer un compromis entre le nombre d'additionneurs et le nombre d'étapes additionnelles, c'est-à-dire entre la surface et la vitesse
- Génère un contrôleur satisfaisant la contrainte de retard de sorte que le nombre d'additionneurs / soustracteurs soit minimisé ; etc.

l'arithmétique radix- $2^r$  est un outil mathématique simple et puissant qui pourrait être exploré pour la résolution de problème de la multiplication par une constante ,les chapitres qui suivent montrer l'efficacité de cet algorithme.

## II.5 Conclusion

Dans ce chapitre nous avons introduit le problème de la multiplication par une constante simple et multiple MCM /SCM et nous avons présenté brièvement la nouvelle solution algorithmique radix  $2^r$  pour la résolution de ce problème

Le chapitre qui suit entame la multiplication MCM/SCM au niveau bloc d'additionneur on montre l'efficacité de cette résolution algorithmique radix  $2^r$

# **Chapitre III :**

**SCM/MCM au niveau bloc  
d'additionneur**

### III.1 Introduction

Contrairement aux heuristiques qui font la résolution des problèmes MCM/SCM, RADIX- $2^r$  offre une réduction très importante au nombre d'additions nécessaire pour représenter un problème SCM/MCM. Cette optimisation influe directement sur l'architecture d'un système qui est basé sur la multiplication par multiple constantes concernant la surface, la vitesse, le chemin critique, ainsi que la consommation d'énergie.

Dans ce chapitre on va traiter l'explication du principe de l'algorithme RADIX- $2^r$  pour la multiplication par constantes, puis la solution SCM/MCM développé au niveau bloc d'additionneur.

### III.2 RADIX- $2^r$ pour la multiplication par une constante SCM/MCM

L'arithmétique RADIX- $2^r$  est une heuristique développée pour minimiser le nombre des additions dans la multiplication par une constante. Sa complexité est sous linéaire [22]. Elle permet de générer un recodage d'une constante  $C$  afin d'avoir un produit  $C \times X$  qui utilise un nombre d'additions optimisé.

La représentation RADIX- $2^r$  est développée par SAM en 1990 [8]. Dans cette technique un nombre représenté en complément à 2 de  $N$ -bits est écrit comme suit:

$$C = \sum_{j=0}^{(N+1)/r-1} (c_{rj-1} + 2^0 c_{rj} + 2^1 c_{rj+1} + 2^2 c_{rj+2} + \dots + 2^{r-2} c_{rj+r-2} - 2^{r-1} c_{rj+r-1}) \times 2^{rj}$$

$$C = \sum_{j=0}^{(N+1)/r-1} Q_j \times 2^{rj} \quad (\text{III.1})$$

Où  $x_{-1} = 0$  et  $r \in \mathbb{N}$

A des fins de simplicité et sans perte de généralité, on considère que  $r$  est un diviseur de  $N$  eq (III.1), la représentation complément à 2 de  $x$  est divisée en  $n+1/r$  complément à 2 segments ( $Q_j$ ), représentés sur  $r$  bits car il va de  $2^0$  jusqu'à  $2^{r-1}$ . Toutefois  $Q_j$  a besoin d'un bit supplémentaire ( $c_{rj-1}$ ) égal au bit de poids fort du segment précédent ( $Q_{j-1}$ ), et donc chaque 2 segments adjacents ont un bit en commun. L'ensemble des chiffres  $DS(2^r)$  correspond à (III.1) tel que :

$$Q_j \in DS(2^r) = \{-2^{r-1}, -2^{r-1} + 1, \dots, 2^{r-1} - 1, 2^{r-1}\}$$

Le produit devient

$$C \times X = \sum_{j=0}^{(N+1)/r-1} X \times Q_j \times 2^{rj} \quad (\text{III.2})$$

Le signe des termes  $Q_j$  correspond au bit  $C_{rj+r-1}$  et  $|Q_j| = 2^{kj} \times m_j$  avec  $kj \in \{1, 2, 3, \dots, r-1\}$  et  $m_j \in OM\{2^r\} \cup \{0, 1\}$  où  $OM(2^r) = \{3, 5, 7, \dots, 2^{r-1}-1\}$ , dans RADIX- $2^r$   $OM(2^r)$  est l'ensemble des nombres positifs impaires avec  $|OM(2^r)| = 2^{r-2}-1$

Chaque  $Q_j$  comprend  $r+1$  bits. Le nombre total des différentes combinaisons des bits est  $2^{r+1}$ , à partir de (III.1) seulement deux combinaisons produisent  $Q_j=0$ , en cas où tous les  $r+1$  bits sont à '0' ou '1'

Finalement on peut exprimer le produit comme suit :

$$C \times X = \sum_{j=0}^{(N+1)/r-1} (-1)^{C_{rj+r-1}} \times (m_j \times X) \times 2^{rj+kj} \quad (\text{III.3})$$

Contrairement à la multiplication par une variable ( $Y \times X$ ) où tous les produits partiels ( $m_j \times X$ ) doivent être précalculés, dans la multiplication par une constante ( $C \times X$ ) seulement les sous-ensembles sont nécessaires. En réalité, le nombre des produits partiels est le nombre de différentes valeurs de  $m_j$  générées par le processus de décodage des  $(N+1)/r$  segments (les termes  $Q_j$ ). Donc la génération des produits partiels (PP) consiste premièrement si  $m_j \neq 0$ , lors du calcul des PP  $m_j \times X$  s'il n'est pas précalculer avant, ensuite il est soumis à un décalage à gauche matériel de  $r_j+k_j$  positions. Finalement la négation de  $(-1)^{C_{rj+r-1}}$  dépend du bit  $c_{rj+r-1}$ .

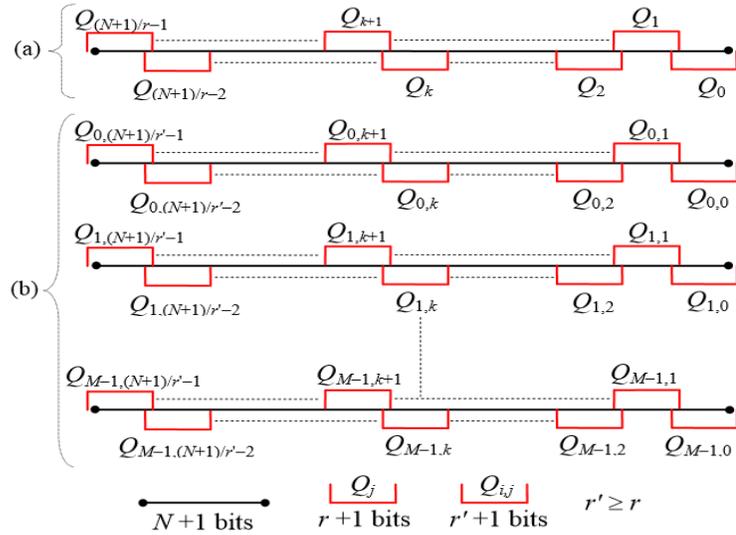


Figure III.1 Partitionnement des segments d'une constante de n bits sous RADIX-2<sup>r</sup>[25]

### III.2.1 Nombre maximal d'additions pour une constante de N bits

D'un côté il y a  $(N+1)/r$  itérations dans (III.3), chaque itération génère un produit partiel, donc le nombre maximal de PP est  $(N+1)/r$  qui requiert au max  $N_{pp} = \frac{(N+1)}{r} - 1$  additions. D'un autre côté, un maximum de  $2^{r-1}-1$  PP non triviaux  $\{3.X, 5.X, \dots, (2^{r-1}-1).X\}$  peut être invoqué lors du processus de génération des PP. Ils sont construits en utilisant la méthode binaire, du bit de poids le plus faible au le bit de poids le plus fort, et pour cela, les éléments de  $m_j$   $3, 5, \dots, 2^{r-1}-1$  sont construits un par l'autre, chaque fois en utilisant une seule addition entre les éléments qui sont déjà construits avec une puissance de 2.

Le processus de construction des PP non triviaux est illustré par l'exemple suivant ou en prend RADIX-2<sup>6</sup> :

$$\begin{aligned} \text{OM}(26) &= \{ 1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31 \} \\ &= \{1\} \cup \{2^1 + 1 = 3\} \cup \{2^2 + 1 = 5, 2^2 + 3 = 7\} \cup \{2^3 + 1 = 9, 2^3 + 3 = 11, \\ &\quad 2^3 + 5 = 13, 2^3 + 7 = 15\} \cup \{2^4 + 1 = 17, 2^4 + 3 = 19, 2^4 + 5 = 21, \\ &\quad 2^4 + 7 = 23, 2^4 + 9 = 25, 2^4 + 11 = 27, 2^4 + 13 = 29, 2^4 + 15 = 31\} . \end{aligned}$$

Alors, les PP ( $m_j \times X$ ) correspondant à OM (2<sup>6</sup>) sont calculés ultérieurement dans l'ordre suivant (6-2=4 étapes) :

$$\{3 \times X\}; \{5 \times X, 7 \times X\}; \{9 \times X, 11 \times X, 13 \times X, 15 \times X\}; \{17 \times X, 19 \times X, 21 \times X, 23 \times X, 25 \times X, 27 \times X, 29 \times X, 31 \times X\}.$$

Donc il est clair que l'EQ (III.3) contient que l'opération d'addition, de soustraction, et de décalage à gauche. La figure III.2 représente tous les détails nécessaires à l'implémentation hardware

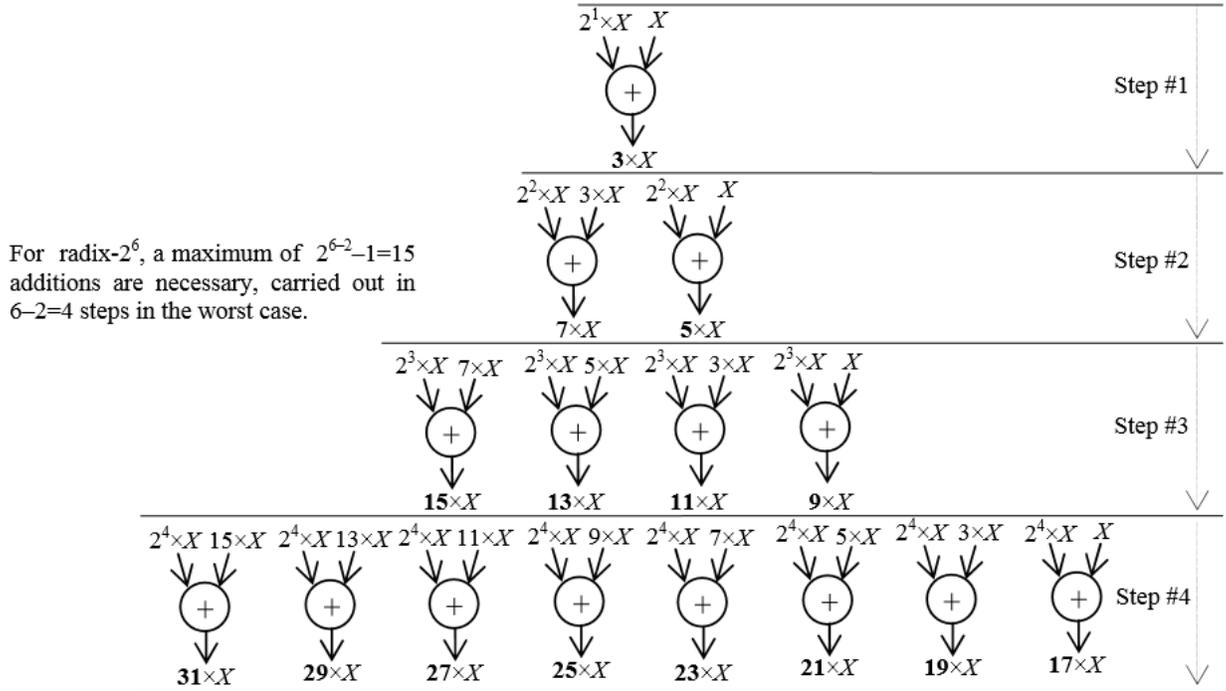


Figure III.2 Ordonnancement séquentiel de l'ensemble des produit partiels nécessaire pour RADIX-2<sup>r</sup>[22]

Par conséquent, le nombre total d'additions requis par Radix-2<sup>r</sup> est égal à :

$$Upb(r) = M_{om} + M_{pp} = \left\lceil \frac{N+1}{r} + 2^{r-2} - 2 \right\rceil \tag{III.4}$$

$Upb(r)$  est minimal pour  $r = 2.W(\sqrt{(N+1).log(2)})/log(2)$  où  $W$  est la fonction de Lambert qu'est la réciproque de la fonction de variable complexe  $f$  définie par  $f(w) = we^w$ , c'est-à-dire que pour tous nombres complexes  $z$  et  $w$ , nous avons :

$$z = we^w \Leftrightarrow w = W(z)$$

Le minimum est obtenu pour un des deux nombres entiers les plus proches de  $r$  et les deux nombres doivent être testés. La figure III.3 représente les limites supérieures "Upper-bounds" des nombres d'additions pour CSD, DBNS, et RADIX-2<sup>r</sup>.

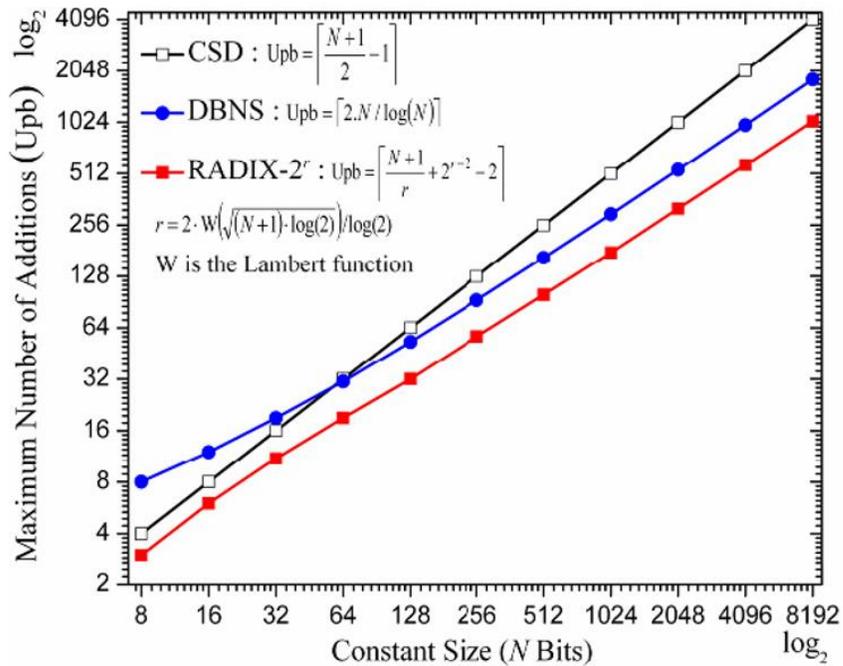


Figure III.3 Comparaison des limites supérieurs pour une constante de N bits

En ce qui concerne le nombre moyen d'additions, il a été calculé de manière exhaustive pour des valeurs de C variant de 0 à  $2^N-1$ , pour  $N = 8, 16, 24$  et  $32$ . Mais pour  $N = 64$ , nous avons calculé la moyenne avec  $10^5, 10^6, 10^9$  et  $10^{10}$  valeurs aléatoires de C uniformément distribuées. Alors que la différence entre les quatre résultats obtenus sont insignifiants ( $<10^{-3}$ ), la valeur Moyenne est autour de 15.7165 additions [22]. Les résultats sont présentés dans le tableau III.1.

Pour  $N = 64$ , RADIX-2r utilise 23,12% en moyennes moins d'additions que CSD. Ce gain semble progresser linéairement pour des valeurs faibles de N.

Tableau III.1 Nombre moyen d'additions de RADIX-2r et CSD[22]

Longueur de bit N de constante	CSD		RADIX-2 <sup>r</sup>		Economie (Avg,%)
	Avg	Upb	Avg	Upd	
8	1.7882	4	1.8645	3	-4.2668 <sup>+</sup>
16	4.4445	8	4.5127	6	-1.5344 <sup>+</sup>
24	7.1111	12	6.7994	9	4.3832
32	9.7777	16	8.9627	11	8.3352
64	20.4444	32	15.7165*	19	23.1256

La figure III.4 représente les valeurs moyennes des différentes heuristiques

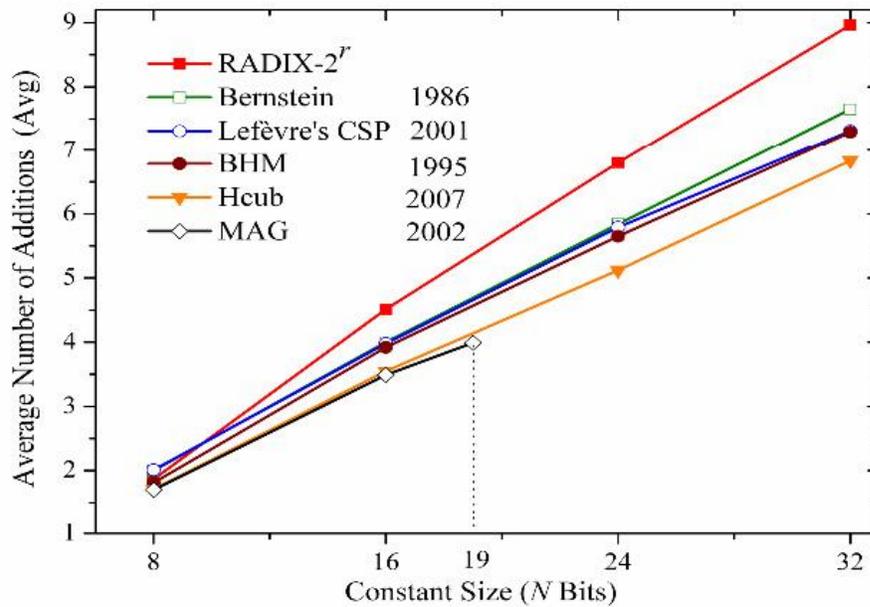


Figure III.4 Longueur de bits N d'une constante[22].

Le tableau III.2 contient les différentes relations de RADIX-2<sup>r</sup> pour un nombre M de constantes 5MCM) a longueur de bits variable.

Tableau III.2 Equations de RADIX-2r pour des constantes non négatives de différents longueur de bit[25]

	Equations
Cout d'additionneurs	$Upb(r') = \sum_{i=0}^{M-1} \left\lceil \frac{(N+1)}{r'} \right\rceil + 2^{r'-2} - 1 - M$ , avec r'=r1 ou r'=r2
Profondeur d'additionneurs	$Ath = \left\lceil \frac{\max(N_i) + 1}{r'} \right\rceil + r' - 3$ , avec i=0, ..., M-1 et r'=r1 ou r'=r2
La moyenne	$-M + Avg_{pp} + Avg_{om} \leq Avg(r') \leq -M - 1 + Avg_{pp} + 2^{r'-2}$ avec

	$Avg_{pp} = (1 - 2^{-r'}) \times \sum_{i=0}^{M-1} \left\lceil \frac{N_i + 1}{r'} \right\rceil$ $Avg_{om} = \sum_{j=0}^{\sum_{i=0}^{M-1} \left\lceil \frac{N_i + 1}{r'} \right\rceil - 1} \left\{ \sum_{k=1}^{2^{r'-2}-1} P(m_{jk}) \times [1 - P(m_{jk})]^j \right\}, \text{ avec}$ $P(m_{jk}) = \frac{\log_2 \left\lceil \frac{2^{r'-1}}{2 \times k + 1} \right\rceil}{2^{r'-1}} \text{ and } r'=r1 \text{ ou } r'=r2$
$r1 = 2.W \left[ \sqrt{\left( \sum_{i=0}^{M-1} (N_i + 1) \right) \cdot \log(2)} \right] / \log(2)$	$r2 = W \left[ 4 \cdot \sqrt{\left( \sum_{i=0}^{M-1} (N_i + 1) \right) \cdot \log(2)} \right] / \log(2)$

Contrairement aux autres algorithmes MCM, chaque constante dans RADIX  $2^r$  est implémentée à part, indépendamment des autres, mais tous les autres constantes utilisent les même PP.

### III.2.2 Les caractéristiques de RADIX- $2^r$

#### III.2.2.1 Totalement prévisible

La limite supérieure (Upb) en nombre d'additions, de profondeur de l'additionneur (Ath) et la moyenne (Avg), sont connus par des formules analytiques exactes. Cette caractéristique non seulement permet aux concepteurs d'avoir une image avant la mise en œuvre sur la surface, la vitesse, et la consommation de puissance, mais elle permet également aux outils de synthèse de satisfaire rapidement les contraintes d'utilisateur et d'effectuer les compromis nécessaires sans les rétroactions « sans fin » à la recherche de la solution appropriée. Notez qu'aucun des algorithmes de MCM existants est prévisible en Upb, Ath ou Avg [32].

#### III.2.2.2 Sous-linéaire

Pour un nombre de constantes non négatives donné M avec une longueur de bit N, RADIX- $2r$  présente une complexité sous-linéaire  $O(M \times N / r)$ , où r est une fonction de (M, N). Cela signifie qu'il n'a aucune limitation concernant le couple (M, N). Pour les problèmes de grande complexité ( $M \times N \gg$ ), RADIX- $2r$  est très probablement le seul qui soit à même capable de fonctionner. Notez que les meilleures heuristiques MCM ont une complexité polynomiale ou exponentielle[32].

### III.2.2.3 Solution haute vitesse et faible consommation de puissance

Profondeur d'additionneur(Ath) n'est pas seulement une mesure de chemin critique (vitesse), mais aussi un bon indicateur de la consommation d'énergie.

Les algorithmes de MCM existants ne peuvent pas rivaliser avec Radix-2r, car il permet une réduction logarithmique de Ath, alors que c'est impossible dans les autres algorithmes en raison des additions partagés [32].

### III.3 Les spécifications d'entrée du système par l'application RADIX-2<sup>r</sup>

La solution RADIX-2r est une application, développée au CDTA, et est basée sur l'arithmétique de multiplication par multiple constantes de Ms A.K Oudjida. Elle prend une matrice de constantes stockées dans un fichier text comme entrée puis elle génère dans un fichier le recodage des constantes du système avec ses spécifications comme sortie.

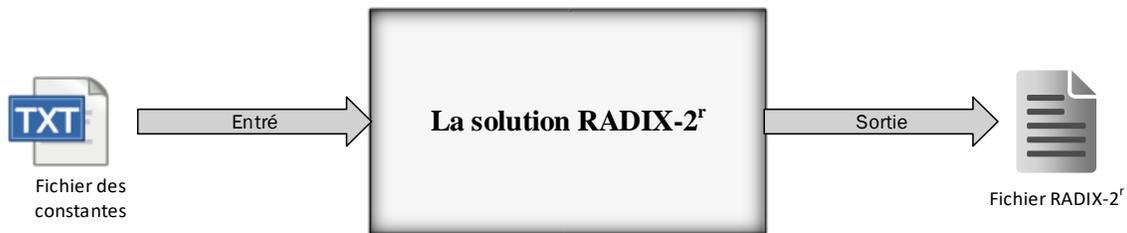


Figure III.5 Le programme RADIX-2r

Cette application donne les résultats suivants :

- Les multiples impaires utilisés ainsi que le recodage des constantes selon l'arithmétique RADIX-2<sup>r</sup> .
- La limite inférieure en nombre d'additionneurs .
- La limite supérieure en cas où les additionneurs sont connectés en série ou en structure d'arbre.
- La longueur de bits maximale et minimale des constantes.
- Le nombre minimal de constantes (Hmin Set) utilisées pour la représentation du SCM/MCM (les odd multiples + les constantes) .
- La valeur optimale de r pour avoir une limite supérieure minimale.
- La valeur de r pour avoir un nombre minimal d'additionneurs.

Le fichier généré par le programme RADIX-2<sup>r</sup> qui contient le recodage ainsi que les spécifications d'un exemple où on a  $C = 10599$  est illustré par la figure III.6

```

16 10599
17
18 *** Radix-2^r MCM solution of the reduced set (Hmin set) including positive/odd constants only ***
19
20 3 = +1<<1 +1<<0
21 7 = +1<<3 -1<<0
22 10599 = +7<<0 +3<<5 -7<<8 +3<<12
23
24 *** Radix-2^r MCM Solution for the user given constants (H set) ***
25
26 10599 = +7<<0 +3<<5 -7<<8 +3<<12
27
28 With:
29
30 7 = +1<<3 -1<<0
31 3 = +1<<1 +1<<0
32
33
34 *** Radix-2^r solution summary ***
35
36 Lower-bound in adder-cost = 3
37 Upper-bound in Radix-2^r = 5
38 Adder-cost in Radix-2^r = 5
39
40 Lower-bound in adder-depth = 3
41 Upper-bound in Radix-2^r (adders connected in series) = 4
42 Upper-bound in Radix-2^r (adders connected in tree structure) = 3
43
44 Maximum constant bit-size (Nmax) in the Hmin set = 14
45 Minimum constant bit-size (Nmin) in the Hmin set = 14
46 Number of constants (M_Hmin) of the Hmin set = 1
47 Optimal value of r (Radix-2^r) for upper-bound = 3
48 Optimal value of r (Radix-2^r) for adder-cost = 4
49

```

Figure III.6 Le recodage de 10599 et sa spécification par l'algorithme RADIX-2'

Cette application permet aussi de générer des filtres FIR d'un MCM donnée en langage verilog au niveau des blocs d'additionneurs.

## III.4 représentation de la solution développée au niveau bloc d'additionneurs

### III.4.1 Les langages utilisés

#### - Langage C

Le langage C a été créé en 1972 par Denis Ritchie avec un objectif relativement limité: écrire un système d'exploitation (UNIX). Mais ses qualités opérationnelles l'ont très vite fait adopter par une large communauté de programmeurs[5].

#### - BASH

Est l'acronyme de *Bourne Again Shell*, est un interpréteur en ligne de commande de type script. Revendiquant ainsi sa filiation forte au shell sh original. D'un autre côté, Bash a aussi intégré de nombreuses fonctionnalités provenant du shell Korn, et même de Tcsh. Il est ainsi devenu un outil très puissant, que la plupart des utilisateurs de Linux emploient comme shell de connexion, on utilise ce langage pour l'ordonnancement des processus de vérification et simulation de système SCM/MCM.

- **VHDL**

Est un langage de description matériel utilisé pour décrire le comportement et la structure des systèmes numériques, l'acronyme VHDL signifie « VHSIC Hardware Description Language », et VHSIC lui-même signifie « Very High Speed Integrated Circuit », toutefois, VHDL est un langage de description matériel à usage général qui peut être utilisé pour décrire et simulé une grande variété de systèmes numériques [7].

### **III.4.2 La description logicielle**

- **Système d'exploitation linux distribution Ubuntu 16.04**

Ubuntu est un système d'exploitation libre, gratuit, sécurisé et convivial, qui peut aisément remplacer ou cohabiter avec les systèmes actuels (Windows, MacOS, GNU/Linux.). À l'aide d'Ubuntu, on peut naviguer sur Internet, lire et écrire des courriels, créer des documents, des présentations et des feuilles de calculs, gérer la bibliothèque multimédia et bien plus encore [28].

- **Editeur de texte G-edit**

Gedit est l'éditeur de texte officiel de l'environnement graphique GNOME Shell, et Unity. Ce logiciel propose une interface simple et facile d'utilisation, développée avec l'aide de la bibliothèque GTK [29].

- **Compilateur GCC**

GCC (GNU Compiler Collection) est une suite de logiciels libres de compilation. On l'utilise dans le monde Linux dès que l'on veut transcrire du code source en langage machine, c'est le plus répandu des compilateurs. La suite gère le C et ses dérivés mais aussi le Java ou encore le Fortran [30].

- **Logiciel de simulation modelsim 12.0 et 13.1**

ModelSim est un environnement complet de simulation et débogage pour la conception de circuits digitaux en ASIC et en FPGA. Il supporte plusieurs langages de description, dont le Verilog, le VHDL et le SystemC.

### III.4.3 Principe

La figure III.7 représente un schéma descriptif de la solution proposée, où le système utilise les spécifications d'entrée d'un système proposé. L'application Radix-2<sup>r</sup> traduit les spécifications d'entrée puis génère une description en langage VHDL optimisée équivalente à la sortie du système.



Figure III.7 la solution SCM/MCM développée

### III.4.4 Stratégie de translation des spécifications du système

#### III.4.4.1 Détection de la zone des multiples impaires

L'extraction des multiples impaires est basée sur la lecture des caractères d'après le fichier d'entrée, la détection de la zone de données désirée fait par les chaînes de caractère qui sont au-dessous et au-dessus de la zone, ces chaînes sont prises comme références, le diagramme de la figure III.8 décrit le processus de détection des multiples impaires.

#### III.4.4.2 Translation des caractères

La zone des données détectée est constituée purement par des chaînes de caractères. Les opérations arithmétiques ne peuvent pas effectuer alors l'étape qui suit. Il faut donc la traduction des caractères vers des types de données bien spécifiés. Cette étape se fait par une série de tests pour assurer l'extraction et la conversion de tous les chiffres et les opérations qui créent tous les multiples impaires. Chaque ligne contient un multiple impair et chaque un de ces multiples est calculé selon ces coefficients  $m_j$  et ces facteurs de décalages  $k_j$ .

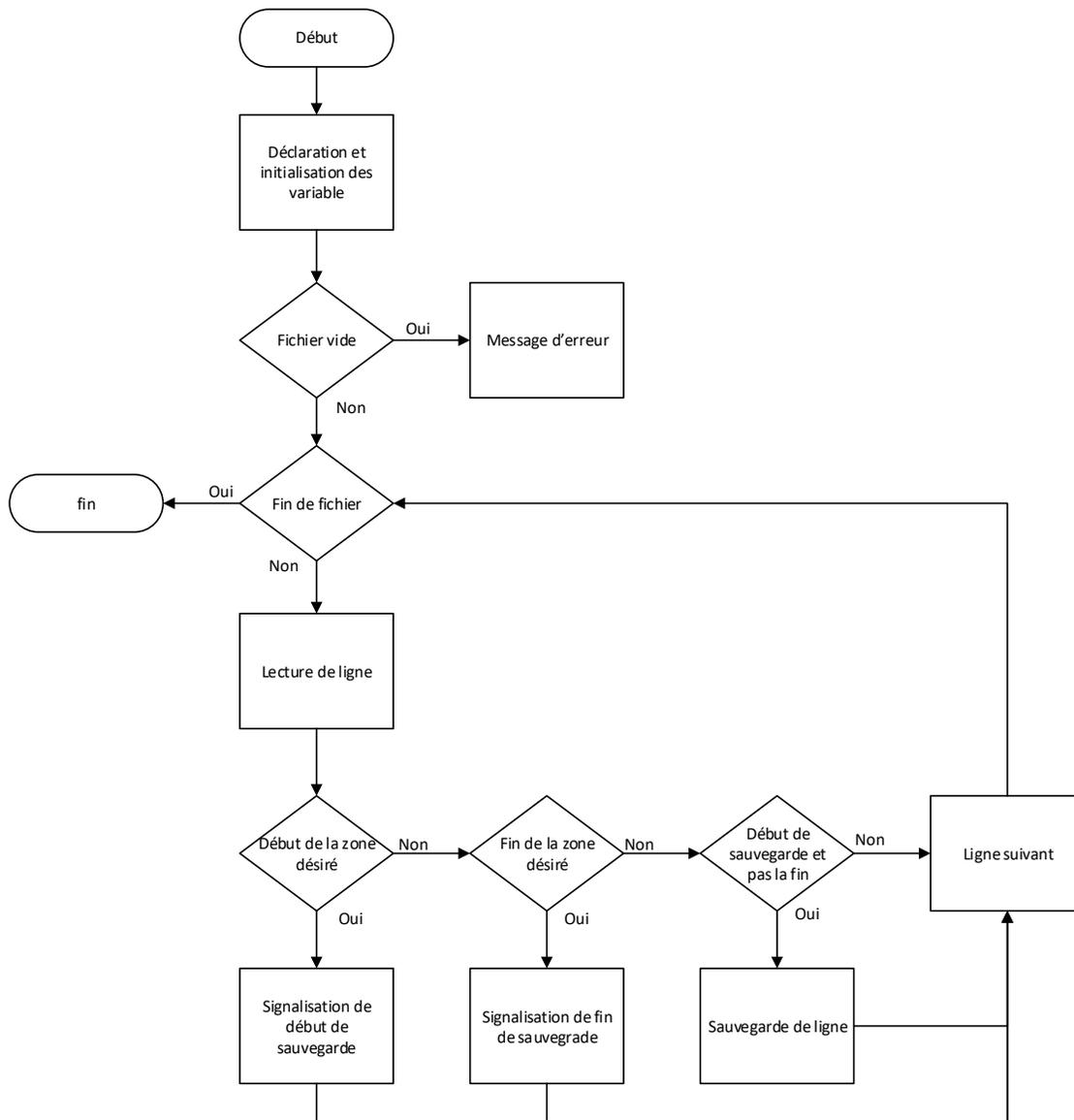


Figure III.8 L'organigramme de la détection de la zone des multiples impaires

### III.4.4.3 Le stockage

Le stockage se fait par l'utilisation de listes chaînées des données. Cette liste est une collection des données ordonnées, de taille arbitraire, dont la représentation en mémoire est une succession de cellules faites d'un contenu et un pointeur vers une autre cellule comme illustré dans la figure III.9.

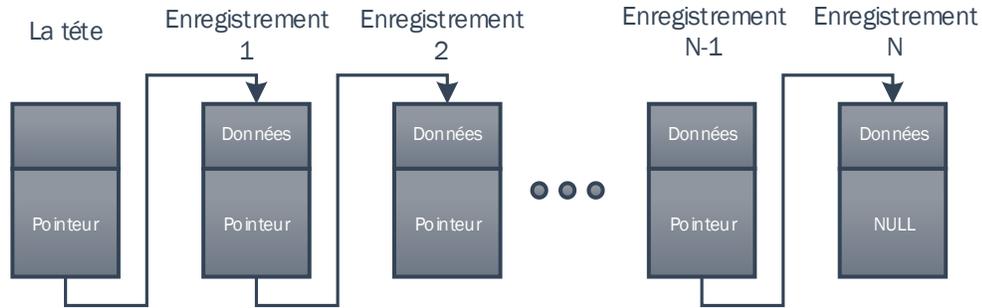


Figure III.9 Le principe de la structure chaine de liste

L'utilisation de ce type de liste permet l'allocation dynamique de mémoire selon le nombre de multiples impairs. Au contraire l'utilisation des enregistrements ou bien des matrices implique une allocation statique ou il exige une taille prédéfinie, le principe de stockage et de conversion des données est présenté sur la figure III.10.

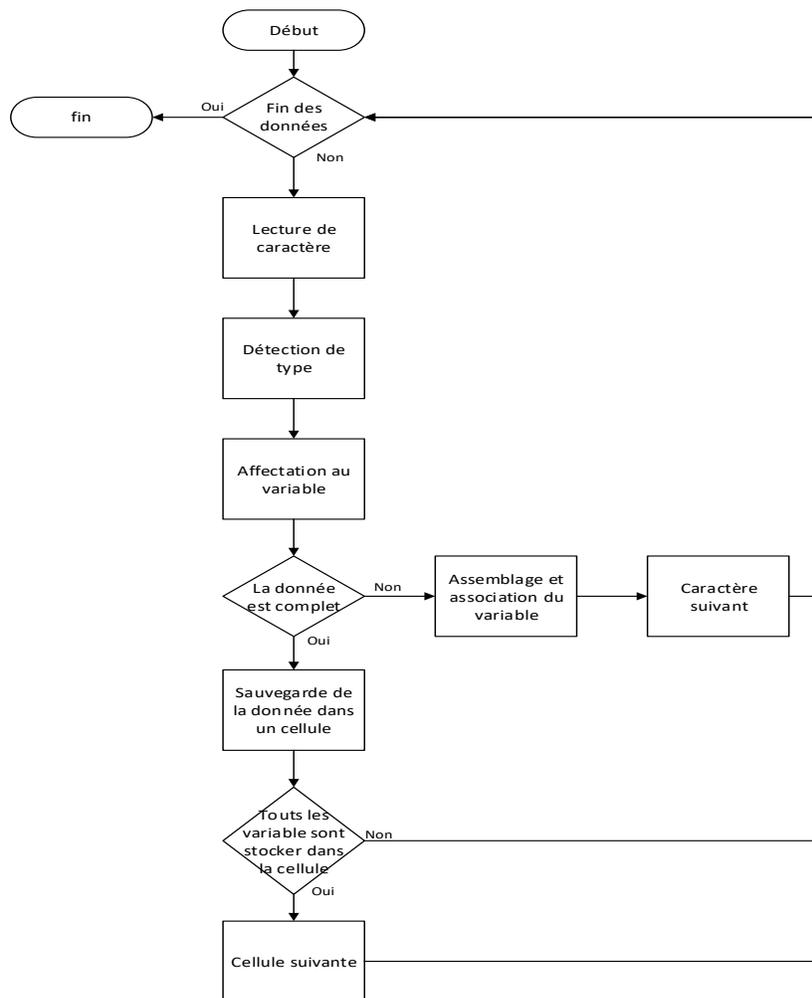


Figure III.10 L'organigramme de conversion et de stockage des coefficients

Après le stockage des  $m_j$  et  $k_j$  de chaque multiple impair ainsi que le calcul de résultat de chaque ligne, le nombre des bits de chaque bloc d'additionneur lui-même doit être calculé et stocké par la formule :

$$bit\_length = Xb + \log_2 C$$

Ce paramètre joue un rôle très important lors de la génération des signaux intermédiaires, la sortie, ainsi que la description de l'architecture de chaque multiple impair.

### III.4.5 Stratégie de génération de la description VHDL

La description matérielle est basée sur la décomposition du système sous des composant élémentaire. Chaque composant représente un bloc d'additionneur d'un multiple impair. Le nombre maximal des composants (bloc d'additionneur) est  $N_{IP} = 2^{r-2} - 1$ . Dans cet exemple  $r=5$  donc le nombre des composant elementaire ne dépasse pas  $2^{5-2} - 1 = 7$ .

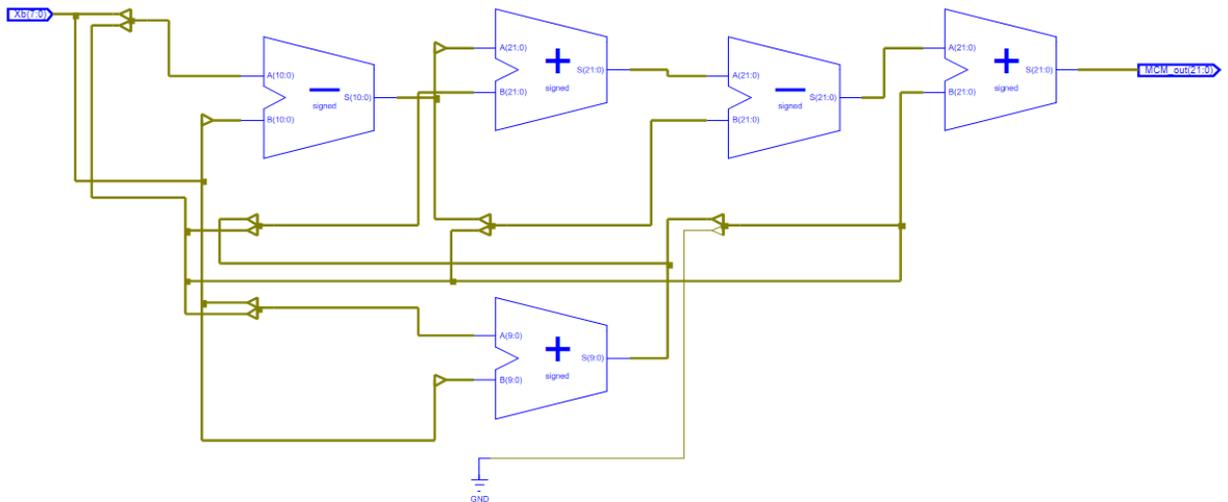


Figure III.11 architecture d'un exemple d'une solution SCM au niveau bloc d'additionneur

Le model VHDL construit est basé sur le type de description flot de donnée « *data flow* », ou on utilise des assignations directes. Tous les instructions sont exécutées en concurrence. On utilise la bibliothèque standard du langage IEEE avec le paquet standard 1164 ainsi que le paquet des nombres signé pour assurer l'addition des nombre positifs et négatifs. Le diagramme de la figure III.12 représente le processus de génération du code VHDL.

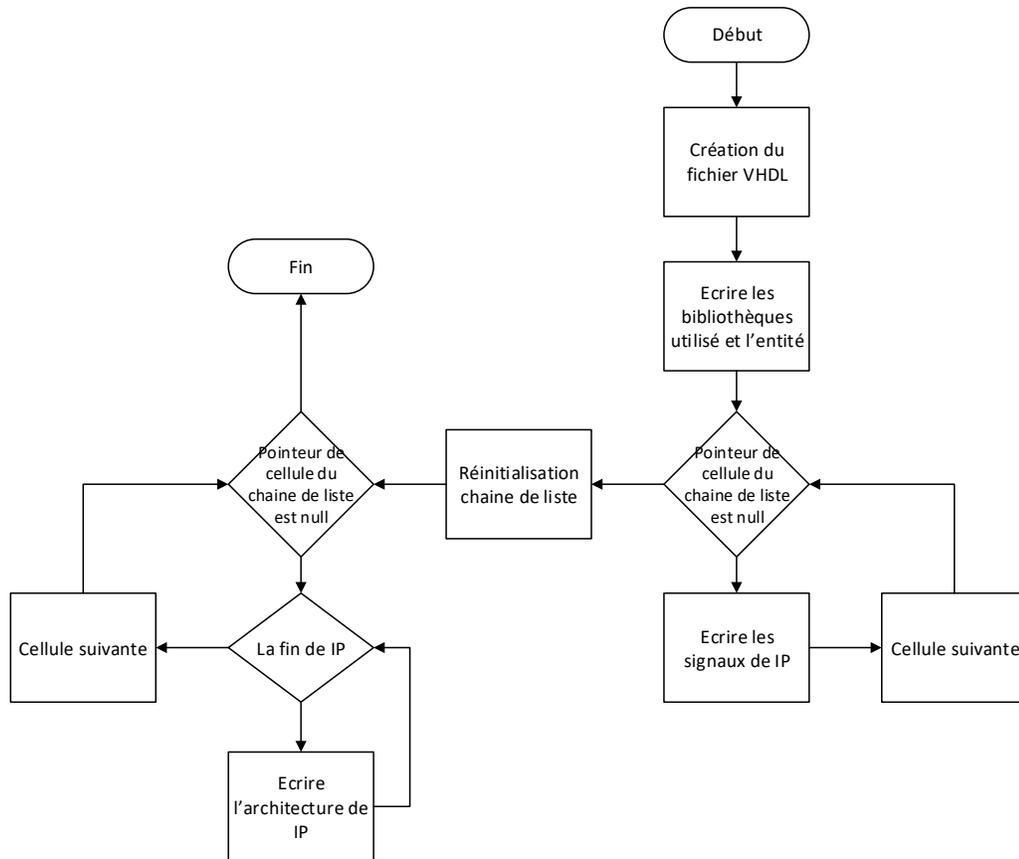


Figure III.12 L'organigramme de génération de la description VHDL

### III.5 Tests et simulations

Afin de s'assurer du bon fonctionnement, on a fait des simulations de 5000 constantes différentes. On a testé toutes les possibilités de l'entrée à l'aide d'un *test bench*. De plus Ce test fait la comparaison du résultat RTL avec la formule mathématique ( $C \times X$ ) comme il est montré au niveau de la figure III.16 (b).

La simulation automatique est faite sous le système d'exploitation linux pour assurer l'ordonancement et faciliter l'organisation des programmes executés . La figure III.13 présente le processus de simulation automatique sous linux.

```

mourad@ubuntu: ~/simulation adder bblock
# -- Loading package TEXTIO
# -- Loading package std_logic_1164
# -- Loading package std_logic_arith
# -- Loading package STD_LOGIC_SIGNED
# -- Compiling entity MCM
# -- Compiling architecture MCM_archi of MCM
# Model Technology ModelSim ALTERA vcom 10.1d Compiler 2012.11 Nov  2 2012
# -- Loading package STANDARD
# -- Loading package TEXTIO
# -- Loading package std_logic_1164
# -- Loading package std_logic_arith
# -- Loading package std_logic_textio
# -- Loading package MATH_REAL
# -- Loading package STD_LOGIC_SIGNED
# -- Loading package NUMERIC_STD
# -- Loading package mypackage
# -- Compiling entity MCM_tb_vhd
# -- Compiling architecture behavior of MCM_tb_vhd
# ** Warning: [4] MCM_tb.vhd(76): (vcom-1207) An abstract literal and an identifier must have a separator between them.
# ** Warning: [4] MCM_tb.vhd(77): (vcom-1207) An abstract literal and an identifier must have a separator between them.
# vsim -lib work -t ips MCM_tb_vhd

```

Figure III.13 Le processus de simulation de constante

**Exemple illustratif :**

On va faire une simulation pour la meme constante utilisé en chapitre 1  $C=10599$ , Le recodage du constante sous le programme RADIX-2<sup>f</sup> est :

$$3 = +1 \ll 1 + 1 \ll 0$$

$$7 = +1 \ll 3 - 1 \ll 0$$

$$10599 = +7 \ll 0 + 3 \ll 5 - 7 \ll 8 + 3 \ll 12$$

La figure III.15 représente une partie de la description VHDL équivalente à cette spécification.

```

20
21 architecture MCM_archi of MCM is
22
23 -- Intermediate signal declaration
24 signal sig_1 : std_logic_vector(7 downto 0);
25 signal sig_3 : std_logic_vector(9 downto 0);
26 signal sig_7 : std_logic_vector(10 downto 0);
27 signal sig_10599 : std_logic_vector(21 downto 0);
28 begin
29 sig_1 <= Xb ;
30 sig_3 <= +(sig_1(7) &
31 sig_1 & '0') +(sig_1(7) & sig_1(7) &
32 sig_1 );
33 sig_7 <= +(
34 sig_1 & "000") -(sig_1(7) & sig_1(7) & sig_1(7) &
35 sig_1 );
36 sig_10599 <= +(sig_7(10) & sig_7(10) & sig_7(10) & sig_7(10) & sig_7(10) & sig_7(10) & sig_7(10)
37 sig_7 ) +(sig_3(9) & sig_3(9) & sig_3(9) & sig_3(9) & sig_3(9) & sig_3(9) & sig_3(9) &
38 sig_3 & "00000") -(sig_7(10) & sig_7(10) & sig_7(10) &
39 sig_7 & "00000000") +(
40 sig_3 & "0000000000000");
41 MCM_out <= Sig_10599 ;
42 end MCM_archi;
43
44 -----end description of mcm solution-----
45

```

Figure III.14 La description VHDL de  $10599 \times X$

Pour la simulation on a utilisé le code qui est déjà donné par la figure III.15, ainsi que le testbench qui génère les stimuli d'entrée ainsi que la comparaison de resultats RTL avec le model mathématique dans le logiciel de simulation modelsim. Les signaux de simulation ainsi que le rapport sont présenté sur la figure III.16 (a) et III.16 (b) respectivement.

Files		Msgs			
/mcm_tb_vhd/Cin	0				
/mcm_tb_vhd/Xb	45	0	1	2	3
/mcm_tb_vhd/MCM...	478953	0	10599	21198	31797

(a)

```

math theory    -74193  =rtl  -74193
math theory    -63594  =rtl  -63594
math theory    -52995  =rtl  -52995
math theory    -42396  =rtl  -42396
math theory    -31797  =rtl  -31797
math theory    -21198  =rtl  -21198
math theory    -10599  =rtl  -10599
simulation succeed|
    
```

(b)

Figure III.15 Le resultat de modelsim et le reeport de simulation de 1059 x X

### III.6 Conclusion

Dans ce chapitre nous avons présenté l'heuristique RADIX-2<sup>r</sup> et ses caractéristiques. On a présenté encore le programme RADIX-2<sup>r</sup> qui permet de calculer et générer la spécification particuliere pour chaque SCM/MCM puis on a presenté la solution que nous avons développée afin de construire la description VHDL équivalente pour chaque cas SCM/MCM.

Le chapitre suivant sera consacré à la présentation de l'heuristique RADIX-2<sup>r</sup> et ses caractéristiques au niveau bit, puis la solution principale développée dans le cadre de notre projet.

# **Chapter IV :**

**SCM/MCM au niveau bit**

### IV.1 Introduction

Les heuristiques au niveau bit produisent de meilleurs résultats d'optimisation en termes de ressources hardware que les heuristiques au niveau bloc d'additionneurs, mais au détriment d'un effort excessif de calcul. Cela les disqualifie de faire face au problème MCM de haute complexité. Un modèle au niveau bit précis est utilisé dans la fonction objective d'un algorithme MCM exact pour estimer la superficie de chaque opération.

Dans ce chapitre nous allons présenter RADIX-2<sup>r</sup> au niveau bit avec son temps d'exécution nécessaire, la partie développement de notre solution optimisée, avec des tests bien simulés qui sont générés d'une manière automatique au niveau bloc d'additionneur et au niveau cellule d'additionneur (1 bit).

### IV.2 Le principe de RADIX-2<sup>r</sup> SCM/MCM au niveau bit

La multiplication par une constante SCM/MCM sous RADIX-2<sup>r</sup> est basée sur l'addition, la soustraction, et le décalage à gauche. Au niveau bloc d'additionneurs (Adder block) l'addition et la soustraction sont faites directement, bien qu'au niveau bit on sait que chaque deux Q<sub>j</sub> adjacent (Q<sub>j</sub>, Q<sub>j+1</sub>) sont décalés un par rapport à l'autre par un décalage à gauche K<sub>j</sub>. Le résultat de l'addition/soustraction donc soit l'addition de tous les bits de Q<sub>j</sub> avec celle de Q<sub>j+1</sub>, alors qu'on peut bénéficier un nombre très important des additionneurs complets par l'élimination des addition/soustraction qui sont pas nécessaires pour l'opération. On ce qui concerne le bit ou il y a le décalage à gauche, dans la majorité des cas on fait le déplacement de ces bits vers le poids faible du résultat plus une addition/soustraction des bits restants.

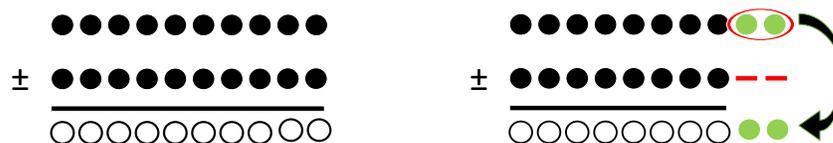


Figure IV.1 le principe d'optimisation au niveau d'additionneurs de 1 bit

#### IV.2.1 Extensions de signe

Nous considérons le cas général de MCM  $\{C_0, C_1, C_2, \dots, C_{M-1}\} \times X$ , Où C<sub>i</sub> est une constante non négative et X représenté au format complément à deux. Cependant, en arithmétique

complément à deux, le signe de tous les opérandes doit être étendu à la longueur en bits du résultat avant toute opération.

Deux produits partiels PPs  $(\dots + 2^{r(j+1)} \times Q_{j+1} \times X + 2^{rj} \times Q_j \times X + \dots)$  exigent une extension de signe SE variant de 1 à  $2r-1$  bits. Quand  $Q_{j+1} < Q_j$  l'extension de signe SE va être minimale et vice versa. Un exemple d'une extension de signe est donné dans la figure 4.2 ou la longueur de bits de X égal a 8,  $r=3$ ,  $j=0$ [24].

$$\begin{cases} Q_j = \pm 1 \Rightarrow 2^{rj} \times Q_j \times X = 2^0 \times 1 \times X = X \\ Q_{j+1} = \pm 2^2 \Rightarrow 2^{r(j+1)} \times Q_{j+1} \times X = 2^3 \times 2^2 \times X = 2^5 \times X \end{cases}$$

$$\begin{cases} Q_j = \pm 2^2 \Rightarrow 2^{rj} \times Q_j \times X = 2^0 \times 2^2 \times X = 2^2 \times X \\ Q_{j+1} = \pm 1 \Rightarrow 2^{r(j+1)} \times Q_{j+1} \times X = 2^3 \times 1 \times X = 2^3 \times X \end{cases}$$

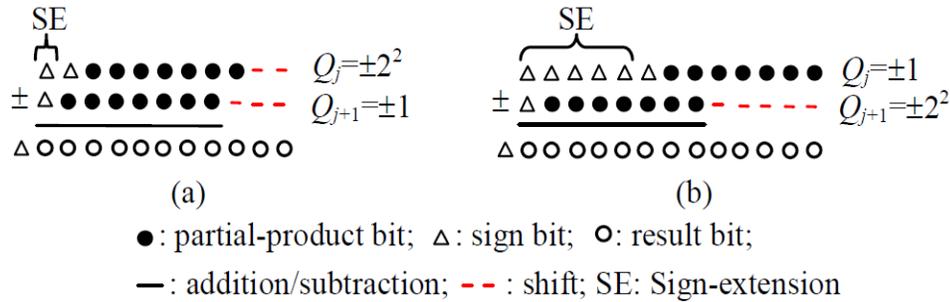


Figure IV.2 la methode d'extension de sign appliqué dans RADIX-2<sup>r</sup>[23]

### IV.2.2 Le nombre total d'additionneurs de 1 bit

Le produit  $C \times X$  requis au max  $X_b + \lceil \log_2 c \rceil$  bits ou  $X_b$  est la longueur de bits de la constante  $C$ . d'un côté les PPs exigent  $\lceil (N+1)/r \rceil - 1$  additions, chaucune consomme au plus  $X_b + r - 1$  additionneurs de type FAs. Alors le nombre max de FAs nécessaire pour la matrice des PPs est  $FA_{pp} = (\lceil (N+1)/r \rceil - 1) \times (X_b + r - 1)$  [11]. (IV.1)

D'un autre côté, les odd multiples consomment au plus

$$FA_{om} = X_b \times (2^{r-1} - 1) + (r-2) \times 2^{r-3} - 1$$
 (IV.2)

Donc C x X a besoin au maximum de FA<sub>pp</sub> + FA<sub>om</sub> d'additionneur de type FA.

**IV.2.2.1 En cas du MCM**

$$FA_{pp} = \sum_{i=0}^{M-1} (\lceil (N+1) / r \rceil - 1) \times (X_b + r - 1) \tag{IV.3}$$

$$FA_{om} = X_b \times (2^{r-1} - 1) + (r - 2) \times 2^{r-3} - 1 \tag{IV.4}$$

$$r = 2.W \left\lceil \sqrt{\left( \sum_{i=0}^{M-1} (N_i + 1) \right) \times \log(2)} \right\rceil / \log(2) \tag{IV.5}$$

**IV.2.3 L'avantage de RADIX-2<sup>r</sup> au niveau bit**

RADIX-2<sup>r</sup> est un algorithme efficace quelque soit la taille en bit (N) de la constante. C'est le seul algorithme qui offre un temps d'exécution sous-linéaire  $O(M \times N / r')$  [25] par rapport à la complexité du problème. Par contre les autres algorithmes offrent des complexités polynomiales ou carrément exponentielles. Quand la constante est très grande, le temps d'exécution sera très important. Le tableau IV.1 représente la complexité des algorithmes MCM.

Tableau IV.1 la complexité de calcul des algorithmes SCM/MCM [25]

<b>Hcub</b>	$O(M^4 \times N^5 \times \log(M \times N) + M^3 \times N^6)$
<b>BHM</b>	$O(M^3 \times N^4)$
<b>Lefèvre CSP</b>	$O(M^3 \times N^3)$
<b>RAG<sub>n</sub></b>	$O(M^2 \times N^3 \times \log(M \times N))$
<b>MAD<sub>c</sub></b>	$O(M \times 2^N)$
<b>NAIAD</b>	$O(M \times 2^N)$
<b>SIREN</b>	$O(M \times 2^N)$
<b>CSD</b>	$O(M \times N)$
<b>R3</b>	$O(M \times N)$
<b>RADIX-2<sup>r</sup></b>	$O(M \times N / r')$

### IV.3 Le developpement d'une solution SCM/MCM au niveau bit

Le principe de la solution au niveau bit est le même qu'au niveau bloc d'additionneurs où on utilise les spécification des constantes du système SCM/MCM générées par RADIX-2<sup>r</sup>, puis extraire les coefficients  $m_j$  et  $k_j$  ainsi que les caractéristiques particulières de chaque constante pour construire une description matérielle en VHDL équivalente aux spécifications utilisées. Le nombre d'additionneurs 1 bit (FA) nécessaires ne dépasse jamais le nombre calculé par la formule. La procédure de détection de la zone des multiples impaires, la conversion des données extraites, ainsi que le stockage des coefficients  $m_j$  et  $k_j$  sont les mêmes. Le calcul du nombre d'additionneurs 1 bits du système est fait par l'utilisation des coefficients  $k_j$  et  $m_j$  et est stocké dans la liste chaînée des données.

#### IV.3.1 Stratégie de calcul du nombre des additionneurs 1 bits

Afin d'obtenir de meilleures performances du model en termes de surface occupée, vitesse d'exécution, ainsi que la consommation d'énergie, l'utilisation d'un nombre d'additionneurs exact est nécessaire pour avoir les résultats désirés. C'est pourquoi la procédure de calcul des additionneurs 1 bit utilisé prend une importance majeure lors du développement de la solution

Le calcul du nombre d'additionneurs est basé sur le calcul de la longueur en bits de chaque couple d'opérandes successives par la relation :

$$bit\_length = m_j + \log_2 C \quad (IV.6)$$

Puis la comparaison de ces deux longueurs et l'élimination des additionneurs qui font l'addition avec des zéros (les additionneurs où il y a un décalage). Le seul cas où l'utilisation d'additionneurs complet sans élimination c'est lorsque  $k_j[i] > k_j[i+1]$  et le deuxième opérande est négatif ( $m_j[i+1] < 0$ ), comme le montre la figure IV.3(d), les autre cas (a), (b), (c) permettent l'annulation des additionneurs 1 bit égaux au coefficient de décalage  $k_j$ .

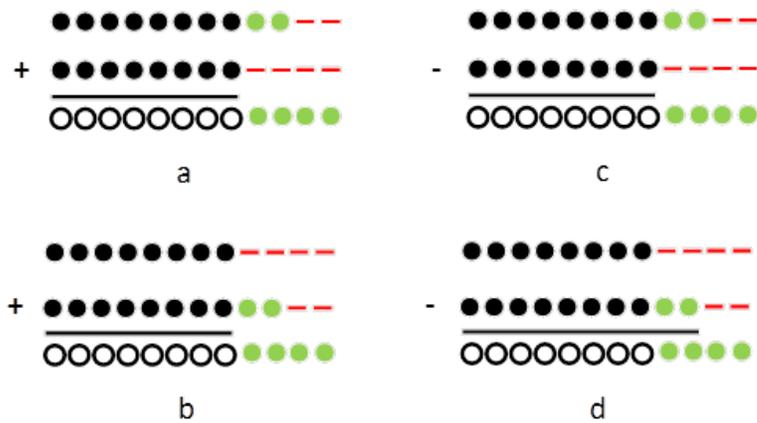


Figure IV.3 les possibilité d'une addition/soustraction en RADIX-2'

**IV.3.1.1 Cas particulier**

Si  $m_j[0] < 0$  on inverse les signes de tous les opérandes (on prend (-) comme facteur)

$$S = -a_0 \ll k_{j_0} + a_1 \ll k_{j_1} - a_2 \ll k_{j_2} \text{ devient } S = a_0 \ll k_{j_0} - a_1 \ll k_{j_1} + a_2 \ll k_{j_2}$$

puis le calcul du nombre d'additionneurs se fait comme les autres cas.

Le diagramme IV.4 représente la technique utilisée pour le calcul des additionneurs.

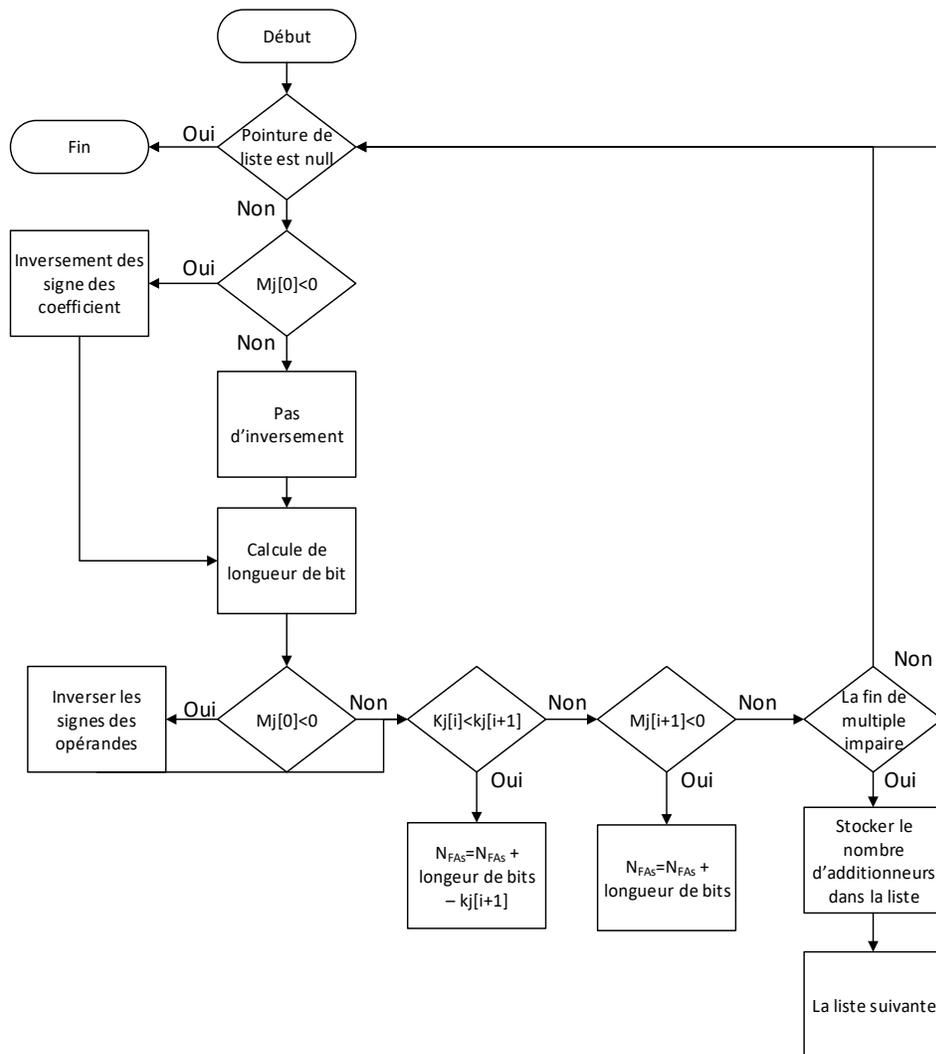


Figure IV.4 l'organigramme de calcul du nombre total d'additionneurs 1 bit

### IV.3.2 La génération du code VHDL

La manipulation au niveau bit consiste à manipuler les données à l'aide des opérations booléennes ET (AND), OU (OR), OU exclusif (XOR) et NON (NOT), ainsi que les décalages logiques et arithmétiques. La génération du code VHDL au niveau bit passe par plusieurs étapes d'analyse et de traitement des cellules de la liste chaînée des données stockées dans le but d'avoir des lignes d'instruction bien précises qui décrivent le système avec une utilisation exact des ressources. On s'appuyant sur une description structurée où on utilise une description d'un composant (additionneur complet 1 bit) puis instancier ce composant afin d'obtenir un bloc d'additionneur optimisé avec une bon fonctionnement de chaque composant élémentaire qui représente une multiplication C x X, le diagramme IV.5 représente les étapes de génération du code VHDL.

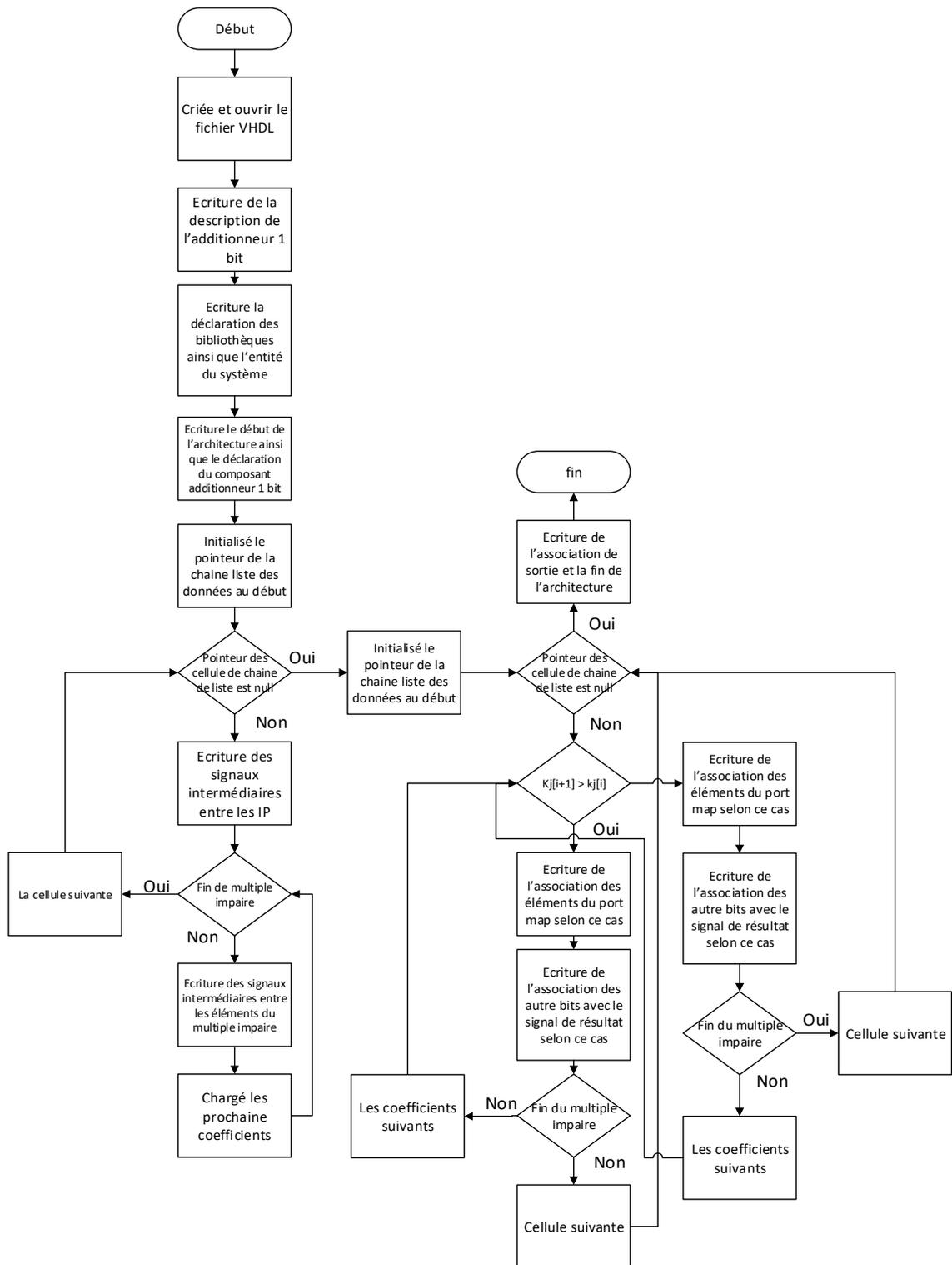


Figure IV.5 organigramme de génération du code VHDL d'un SCM/MCM

## IV.4 La description VHDL du SCM/MCM

### IV.4.1 L'additionneur complet 1 bits (FA)

C'est un circuit qui fait l'addition en bit, il a 3 entrées pour le premier bit, le deuxième bit et le retenu précédent comme un troisième bit. La sortie est constituée de la somme et de la retenue. Le schéma ci-dessus ainsi que le tableau représente le circuit d'un additionneur complet 1 bit et sa table de vérité respectivement.

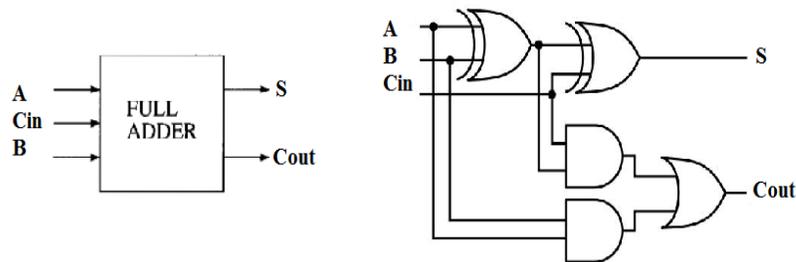


Figure IV.6 un additionneur complet 1 bit

Tableau IV.2 table de vérité d'un additionneur complet 1 bit

A	B	Cin	S	Cout
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

### IV.4.2 L'additionneur soustracteur

L'additionneur/soustracteur est un circuit qui fait les deux opérations (l'addition, la soustraction). La soustraction des deux nombres A et B donnée par  $S = A + B' + 1$ , où B' est le

complément à 1 de B, l'implémentation hardware de cette formule fait par une porte XOR ou l'ensemble des bits de B sera inversé puis à l'aide de l'entrée de la retenue on ajout un 1 pour obtenir un conversion complément à 2. La figure IV.7 représente un exemple d'un additionneur soustracteur de 4 bits. La table de vérité (Tableau IV.3) représente l'inversement des bits lorsque on fait une soustraction (Add/sus =1).

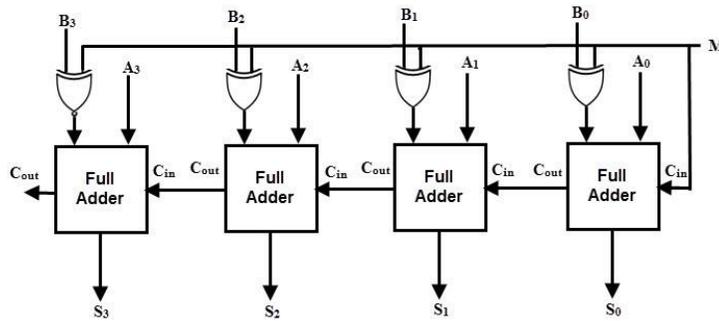


Figure IV.7 le circuit d'un additionneur soustracteur de 4 bits [31]

Tableau IV.3 table de verité d'un inversement par XOR

Add/sus	B	S
0	0	0
0	1	1
1	0	1
1	1	0

### IV.4.3 L'architecture d'un SCM/MCM

Au niveau bit on élimine tous les additionneurs de 1 bit dont nous n'avons pas besoin. Ces additionneurs font la somme des bits avec des zéros, alors qu'il suffit juste de l'assignée à la sortie du signal de résultat. Cette opération ne consomme aucune ressource hardware. La figure IV.8 représente l'architecture d'un exemple d'un multiple impaire au niveau bit par l'utilisation de cette solution.

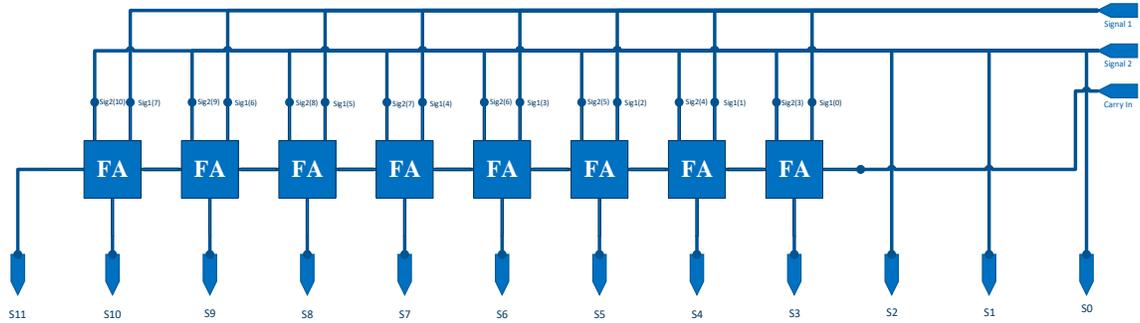


Figure IV.8 L'architecture d'un SCM au niveau bit

On associe les circuits des multiples impaires afin de construire les circuits SCM/MCM mentionnés dans les spécifications. La IV.9 représente un exemple d'une architecture SCM/MCM.

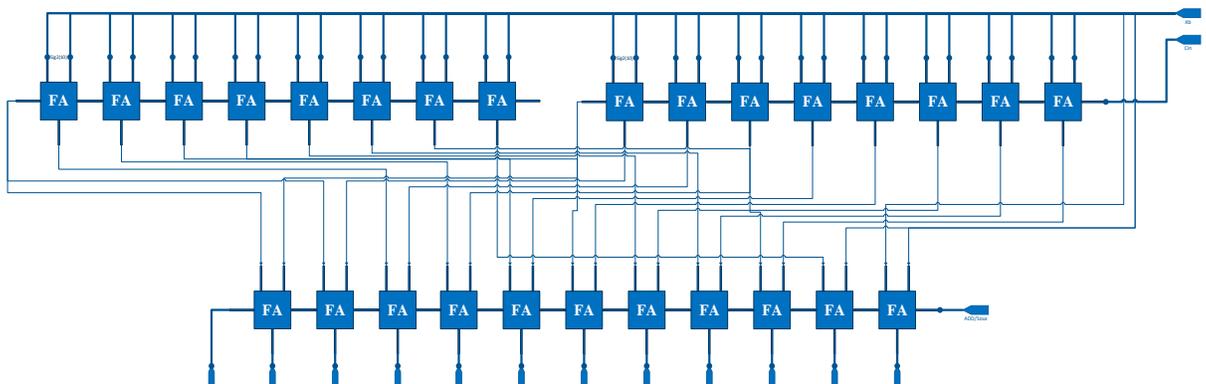


Figure IV.9 l'architecture d'un exemple SCM/MCM

## IV.5 La verification de fonctionnelle

Le but primordial est de prouver que le module va fonctionner dans l'environnement dans lequel il sera intégré. Pour cela, il faut simuler cet environnement, moyennant un *testbench*. Ce dernier permettra de mettre en situation le design, en émulant le mieux possible son environnement. Pour un *testbench* basique, seule l'instanciation du bloc et la génération des stimuli d'entrée sont nécessaires (Figure IV.10).

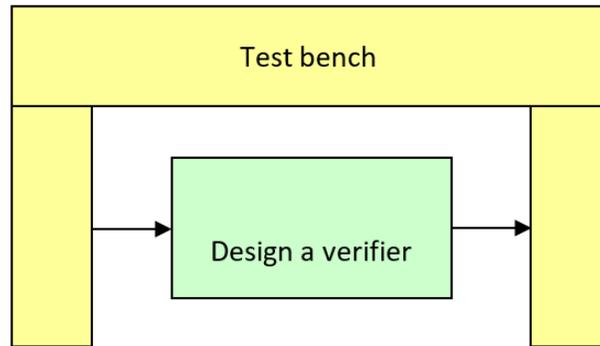


Figure IV.10 test bench

Toutefois, ceci impliquerait une vérification manuelle des résultats, donc une validation incomplète puisqu'une vérification manuelle ne pourra pas couvrir assez de combinaisons d'entrées (des constantes). Les tests effectués à ce niveau ne permettent pas de tester tous les cas possibles d'une part, et d'autre part le concepteur devait lui-même prédire les résultats mathématiques afin de les comparer avec les résultats obtenus au niveau de l'architecture .

Pour cette raison, une vérification fiable ne peut se faire qu'en faisant une comparaison automatique des résultats de sortie théorique et pratique Figure IV.11. Nous pouvons ainsi faire une validation sur un très grand nombre de combinaisons d'entrée, tout en évitant les erreurs humaines.

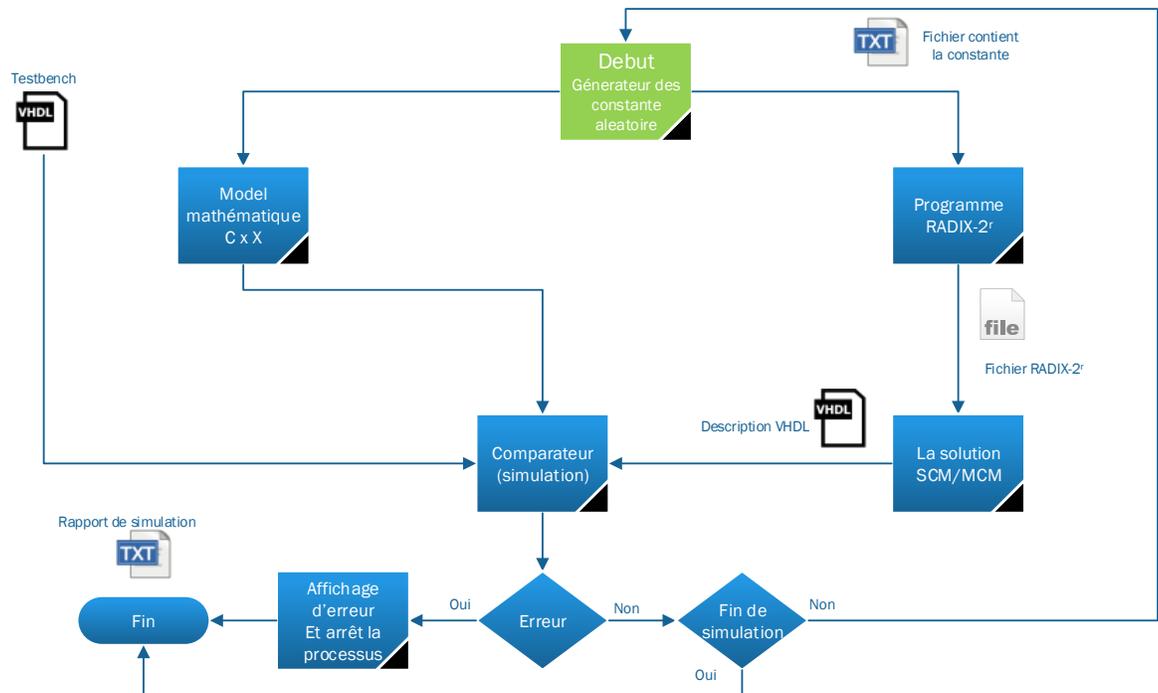


Figure IV.11 principe de la verification fonctionnelle

Le processus de vérification fonctionnelle est ordonné par un script en shell, il commence comme le montre le diagramme, tout d'abord par un programme en C qui génère des constantes aléatoires en sortie. Cette dernière est reliée directement à l'entrée du RADIX-2<sup>r</sup> qui génère lui-même une spécification particulière pour le système C x X désiré. Elle contient le recodage des multiples impairs qui vont être utilisés pour générer la description VHDL équivalente du système par la solution développée.

### IV.5.1 Le script shell

Il est responsable de la vérification fonctionnelle. Il gère les programmes exécutés, ainsi que le lancement du simulateur Modelsim. Le script associe les entrées de chaque programme avec la valeur équivalente. Il génère encore le rapport de simulation qui contient le débogage de la vérification fonctionnelle. En cas d'erreur de simulation ou bien de faux résultats, dans chaque programme le script arrête le processus d'exécution. Les principales parties du script sont illustrées par la figure IV.12.

```

16 do
17
18     ./rand_num # start the constant number generator program
19     if [ -e $constfile ] #verify if the constant file exist and read the constant
20     then
21         echo "loading constant file ..."
22         if [ -s $constfile ]
23         then
24             echo "reading the value ..."
25             read constval < $constfile
26             const=$constval
27             echo $const
28         fi

```

L'exécution de générateur des constantes aléatoire

```

73
74 if [ $equal = false ]
75 then
76     ./radix m constant.txt c | cat # execute the radix program using constant file in the input
77     echo constant.result 8 | ./my_prg # execute developed program using the result file of radix
78     /home/mourad/altera/13.1/modelsim_ase/bin/vsim -c -do ~/simulation/MCM_tb_vhd.fdo # excute modelsim in CLI mode & excute the fdo scripte
79

```

L'exécution de programme RADIX-2<sup>r</sup>, la programme développée, ainsi que le simulateur modelsim

```

111 fi # the succeed condition
112
113 if [ $sim_suc = true ]
114 then
115     echo " #####"
116     echo " # simulation succeed "$counter" are tested #"
117     echo " #####"
118 else
119     echo "there are an error during the simulation please check you log file for more information ."
120     break
121     kill
122
123 fi

```

Figure IV.12 quelques parties de script shell

### IV.5.2 Le Test bench

Le *test bench* est un programme non synthétisable écrit en VHDL ou en Verilog. Il invoque le design puis le simule. Il peut aussi faire des fonctions plus complexes. Il contient une logique pour déterminer les stimulus corrects de design ou comparer les résultats par exemple.

**IV.5.2.1 Les bibliothèques utilisées pour le test bench****Les bibliothèques**

- IEEE
- STD

**Les packages :**

- IEEE.STD\_LOGIC\_1164.ALL (le package standard de VHDL)
- IEEE.STD\_LOGIC\_SIGNED.ALL (la package qui permet la manipulation des nombres signés)
- STD.textio.all (le package qui permet la manipulation de fichiers).
- IEEE.MATH\_REAL.ALL (contient les fonctions mathématiques).

Le *test bench* crée est composé de deux parties : un paquet contient une fonction qui calcule le nombre de bits de sorties pour chaque constante différente (Figure IV.13a), et la deuxième est le corps du *test bench* qui génère les différentes possibilités de l'entrée, puis il compare les résultats RTL (stimulis de sortie) avec le modèle mathématique C x X. Ce processus termine par le débogage des résultats dans un rapport, en cas d'erreur, la simulation sera arrêtée puis il donne dans le rapport un message d'erreur contenant la partie exacte où se trouve l'erreur, les figures IV.13 b et c contiennent les parties principales du *test bench*.

```

13 package body mypackage is
14 function bit_length(const : integer) return integer is
15
16 file const_file : text open read_mode is "constant.txt";
17 variable c : integer ;
18 variable the_length : integer ;
19 variable const_line : line ;
20
21
22 begin
23
24 while not (endfile(const_file)) loop
25   readline (const_file,const_line);
26   read(const_line,c);
27   end loop ;
28
29   the_length := 7 + integer (ceil(log(real (c))/log(real (2)))) ;
30   return the_length ;
31 end bit_length ;
32
33 end mypackage ;
34

```

(a)

```

51 tb : PROCESS
52 BEGIN
53
54
55 -- Wait 100 ns for global reset to finish
56 for i in 0 to 300 loop
57   wait for 1 ns;
58   Xb <= Xb + '1';
59   end loop;
60
61 -- Place stimulus here
62
63   wait; -- Will wait forever
64 END PROCESS tb;
65

```

(b)

```

88 while not (endfile(math_file)) loop
89
90   wait for 1 ns ;
91   readline (math_file,math_line);
92   read(math_line,op);
93
94   math_n := op * c;
95
96   write(logfile_line,string("math theory  "));
97   write (logfile_line,math_n);
98   write (logfile_line,string(" =rt1  "));
99   write (logfile_line,conv_integer(MCM_out));
00   writeline (log_file,logfile_line) ;
01
02   if math_n /= conv_integer(MCM_out) then
03     write(logfile_line,string("error found:  "));
04     write (logfile_line,op);
05     writeline(log_file,logfile_line) ;
06     wait ;
07
08   end if ;
09 end loop ;
10   write (logfile_line,string("the code is succesfully simulated"));
11   writeline (log_file,logfile_line) ;
12   wait ;

```

(c)

Figure IV.13 le code test bench

### IV.5.3 Compilation, élaboration, et simulation

Avant que le model VHDL soit simulé, on doit premièrement faire la compilation figure IV.14, en cas d'erreur la compilation affichera un message d'erreur approprié. Si le code VHDL est conforme à tous les règles, le compilateur génère un code intermédiaire qui peut être utilisé par le simulateur ou le synthétiseur. On utilise le simulateur Modelesim en mode interface ligne de commande CLI. Cette technique offre un contrôle facile et robuste avec une rapidité d'exécution.

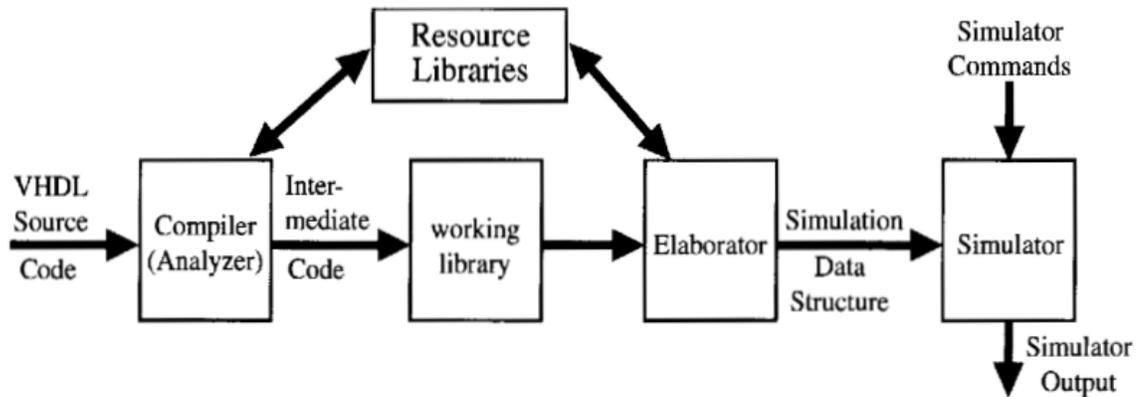


Figure IV.14 compilation, élaboration et simulation du code VHDL [6]

#### IV.5.4 La simulation par modelsim

Les étapes de simulation sont ordonnées dans Modelsim par un script. Les étapes de simulation sont :

- Création d'une bibliothèque pour la simulation
- Compilation des fichiers VHDL (les packages, la description de système, le *test bench*)
- Lancer la simulation par VSIM
- Afficher les signaux
- Déterminer le temps de simulation

Le script utilisé par modelsim est illustré par la figure IV.15.

```

##
vlib work
vcom -explicit -93 "mypackage.vhd"
vcom -explicit -93 "MCM_solution.vhd"
vcom -explicit -93 "MCM_tb.vhd"
vsim -t 1ps -lib work MCM_tb_vhd
view wave
add wave *
do {MCM_tb_vhd.udo}
view structure
view signals
run 1000ns
  
```

Figure IV.15 le script de simulation de modelsim

## IV.6 Resultats et discussion

On a résolu les meme exemples de SCM/MCM au niveau bloc d'additionneurs et au niveau bit (additionneur complet 1 bit) tableau IV.4.

Tableau IV.4 le nombre des additionneurs 1 bit consommés par les deux solutions

MCM/SCM	Nombre Maximum d'additionneurs	Niveau bit	Taux de profit (%)
<b>SCM</b>			
10559 SCM	60	50	16.66
<b>Filtres a faible ordre</b>			
FIR_25_13_12	463	391	15.55
FIR4_30_14_13	586	424	27.64
FIR3_30_14_13	628	469	25.31
FIR1_40_19_12	733	620	15.41
FIR2_40_19_13	775	664	14.32
FIR7_40_19_19	859	642	25.26
<b>Filtre a ordre moyen</b>			
FIR6_60_29_14	1069	841	21.32
FIR8_80_36_14	1069	941	11.97
FIR5_80_39_15	1384	1196	13.58
<b>Filtre a haute ordre</b>			
FIR_279_140_24	5822	4784	17.82
FIR_418_208_22	7156	6128	14.36
FIR_516_256_24	8467	7093	16.22
FIR_631_313_23	9502	7929	16.55
FIR_695_345_24	10399	8665	16.67

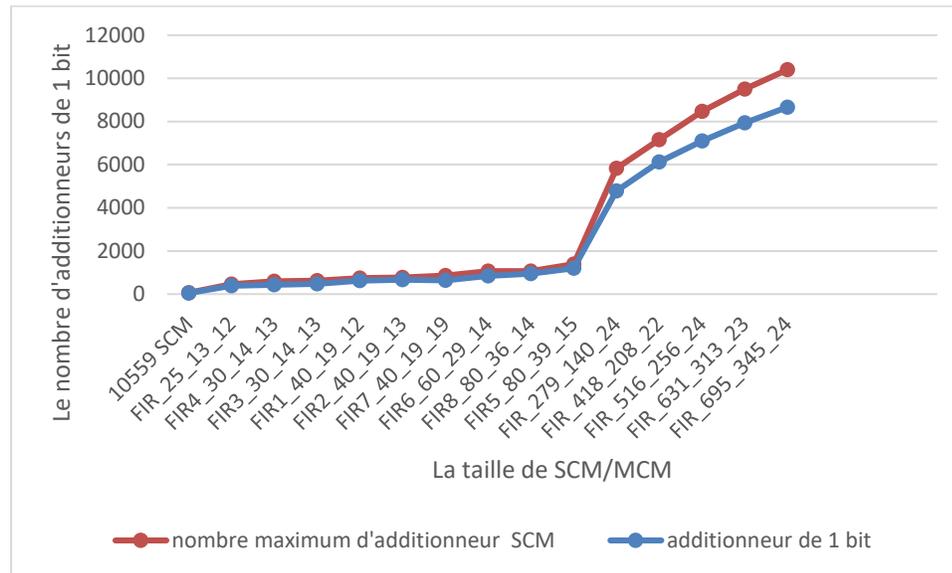


Figure IV.16 Comparaison entre le niveau bloc d'additionneurs et niveau bit

D'après le graphe la réduction du nombre d'additionneurs est plus grande quand l'ordre des filtres ou bien le nombre des MCM sont plus grands. Le graphe de la figure IV.16 représente la variation du rapport niveau bit / niveau bloc. Il montre un rapport compris entre 0.9 et 0.7, et ça offre une réduction de 14% à 27%.

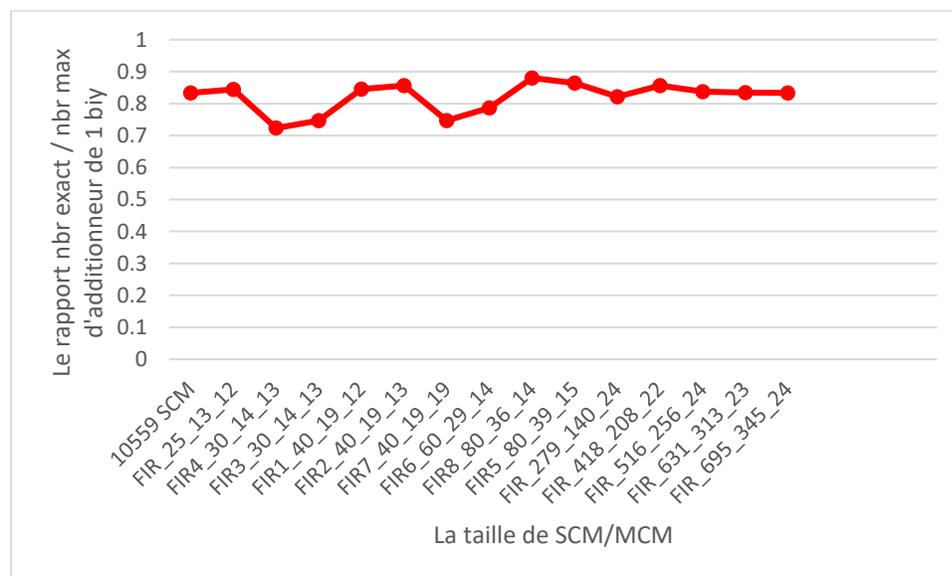


Figure IV.17 Variation d'un nombre d'additionneurs réduits en différents cas

On vu que le niveau bit offre une optimisation très importante, en terme de nombre total d'additionneurs à l'aide d'élimination des additionneurs de 1 bit qui ne sont pas utilisés. Cette optimisation a une influence sur la surface de la puce utilisé (FPGA/ASIC), la vitesse, et la consommation d'énergie.

## **IV.7 Conclusion**

Dans ce chapitre nous avons expliqué l'heuristique de RADIX- $2^r$  au niveau bit, et présenté l'idée de notre solution développée, ainsi que les étapes de conception, avec des tests pour SCM/MCM au niveau bloc d'additionneurs et au niveau bit,

Nous avons constaté que la solution au niveau bit offre une optimisation très importante en termes de ressources hardware qui influence directement la surface, le temps d'exécution, ainsi que l'énergie consommée.

Le chapitre qui suit sera dédié à l'utilisation de cette solution pour construire des filtres numériques RIF optimisé.

# **Chapter V :**

## **Application du SCM/MCM aux filtres RIF**

## V.1 Introduction

Le filtre numérique a une large gamme d'applications, et est un élément clé dans la plupart des systèmes électroniques et plus généralement dans tout dispositif de traitement de signal. La conception matérielle d'un filtre RIF est principalement dominée par des blocs de multiplieurs (SCM/MCM), qui sont implémentés généralement par un réseau d'additionneurs, soustracteurs, et de circuits de décalage (Shifters).

Dans ce chapitre on utilise la solution SCM/MCM au niveau bit développée précédemment pour générer une description matérielle pour un filtre RIF optimisée au niveau de la surface, vitesse, et consommation de puissance.

## V.2 Les filtres RIF

Les filtres à réponse impulsionnelle finie RIF sont des filtres à temps discret avec une durée finie. Ils sont caractérisés par une réponse uniquement basée sur un nombre fini de valeurs de signal d'entrée. Par conséquent, quel que soit le filtre, sa réponse impulsionnelle sera stable et avec une durée finie, dépendant du nombre de coefficients du filtre. Les filtres RIF sont utilisés dans le traitement numérique des signaux [24], Un signal de sortie de N-tap d'un filtre RIF est calculé par:

$$y(n) = \sum_{i=0}^{N-1} h_i \times x(n-i) \quad (\text{V.1})$$

Où N est la longueur du filtre,  $h_i$  sont les coefficients, et  $X_{(n-1)}$  est l'entrée précédente du filtre. La forme directe et transposée d'un filtre RIF sont illustrées par les figures 1(a) et (b), respectivement.

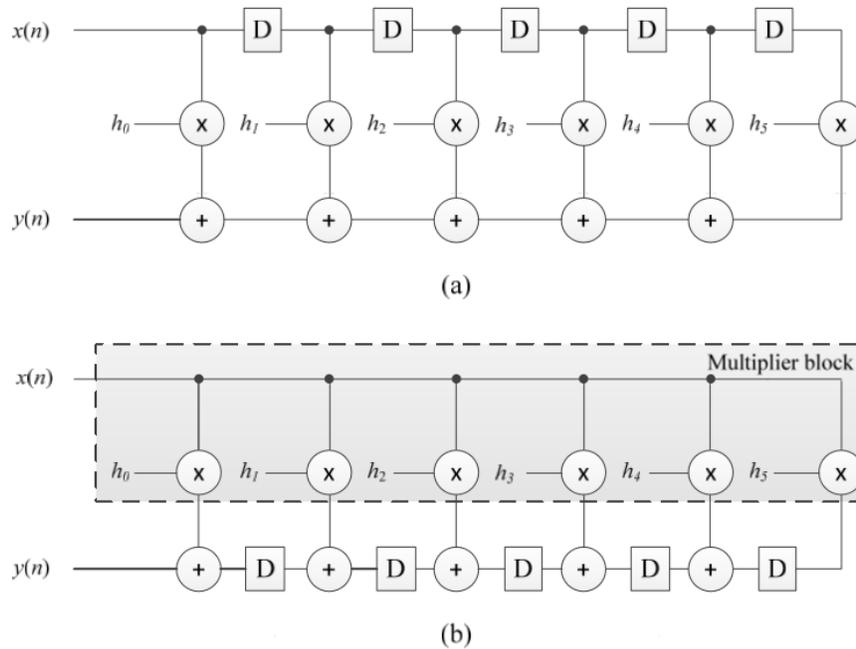


Figure V.1 implementation d'un filtre RIF (a) forme direct (b) forme transposé[24]

Selon la figure V.1 , les deux architectures ont une complexité similaire. Ils ont  $N$  multiplieurs et  $N-1$  bloc d'additionneurs, ainsi que des registres composés principalement par des bascules D.

En général, la réalisation préférée est souvent la forme transposée à cause de sa performance et sa consommation d'énergie qui sont relativement meilleurs. La multiplication des coefficients de filtrage avec l'entrée du filtre à un impact significatif sur la complexité et les performances de la conception car un grand nombre de blocks de multiplication par une constante sont nécessaires [24].

Puisque les coefficients de filtrage sont fixés et déterminés préalablement, la réalisation matérielle de multiplieur est couteuse en terme de surface, délai, et énergie consommé. L'opération de la multiplication par une constante est implémentée généralement de manière sans multiplication (multiplierless) par l'utilisation d'additions, de soustractions, et de circuits de décalage, sachant que l'opération de décalage peut être implémenter juste par des fils, qui n'exigent aucune ressource matérielle (voir le chapitre précédent). Dans ce cas la complexité des filtres RIF dépend uniquement du nombre des additionneurs/soustracteurs utilisés pour implémenter la multiplication par une constante.

### V.3 Les plateformes de réalisation des filtres RIF

La phase de conception d'un filtre numérique consiste à matérialiser l'algorithme de calcul, pour cela, il existe plusieurs possibilités de réalisation [10] :

- Logique câblée (portes logiques, mémoires, ...) par l'utilisation des FPGA/ASIC
- Logique programmée (processeur de traitement du signal "DSP : Digital Signal Processing", microprocesseur "ordinateur").

#### V.3.1 FPGA/ASIC

Un FPGA *Field Programmable Gate Arrays* est un circuit de semi-conducteur basé sur une matrice de blocs élémentaires configurables (CLBs), organisés en lignes et en colonnes et connectés entre eux avec des interconnexions programmables. Les FPGA peuvent être reprogrammables pour une application désirée ou une fonctionnalité requise avant la fabrication. Cette fonctionnalité distingue les FPGA des ASIC.

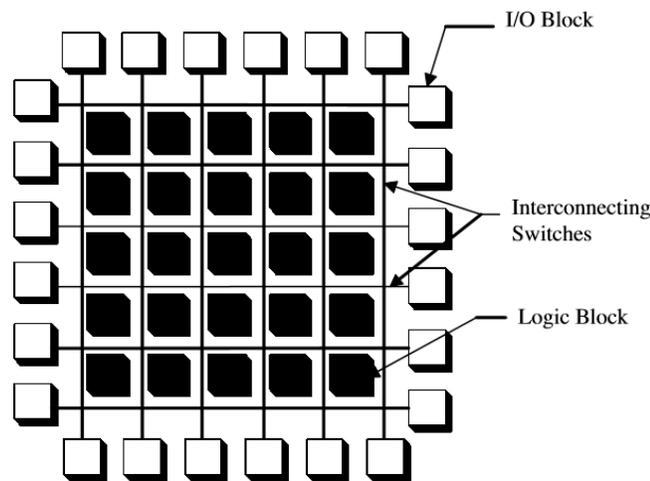


Figure V.2 la structure interne d'une FPGA[7]

ASIC *Application Specific Integrated Circuits* sont des circuits intégrés fabriqués sur mesure pour des tâches de conception spécifiques.

#### V.3.1.1 Caractéristiques des FPGA/ASIC

- Temps de réponse très rapide dans les entrées/sorties
- Des capacités de calcul parallèles massives qui peuvent être exploitées en optimisant l'architecture
- Les multiples cellules MAC peuvent être implémentées sur une seule puce.

### V.3.2 DSP

Un DSP pour Digital Signal Processor en anglais est un processeur spécialisé dans le traitement numérique du signal. Son architecture est optimisée pour traiter une grande quantité de données en parallèle à chaque cycle d'horloge. Ce mode de fonctionnement est très efficace pour traiter des signaux numériques (filtrage, compression, extraction, etc.) .

#### V.3.2.1 Caractéristiques des DSP

- Peu coûteux
- Basés sur l'architecture RISC (Reduced Instruction Set Computer)
- ont des multiplieurs accumulateurs rapides
- Architecture de pipeline multi-stage

Afin de comparer ces deux types de plateforme et pour un même traitement d'image effectué, le tableau ci-dessous illustre les résultats de ce traitement :

Tableau V.1 le traitement d'image dans une FPGA et un DSP

	FPGA	DSP
Le nombre des cycles d'horloge	164354 cycles	10770432 cycles
Temps d'exécution (fréquence d'horloge est 50 Mhz pour FPGA, 600 Mhz pour DSP)	3.2 ms	17.9 ms
Temps d'exécution/pixel	0.19 $\mu$ s	1 $\mu$ s
Utilisation matériels HW	742 éléments logique	67.2 KB
Nombre de ligne de code	132	429
L'énergie consommé	79.97 Mw	540 mW

Il est clair qu'avec une fréquence réduite, une vitesse importante, peu de consommation d'énergie, ainsi qu'un court code on peut implémenter des traitements de signaux sur FPGA contrairement au DSP qui sont relativement plus coûteux en terme d'énergie avec moins de performance. Donc l'utilisation des FPGA/ASIC garantit une haute performance avec un coût en énergie moindre.

## V.4 La structure d'un filtre RIF

La structure d'un filtre RIF se compose d'une chaîne de  $N$  SCM (multiplication par constante),  $N-1$  registres pour les SCM [24], le produit est sommé par une chaîne de  $N-1$  additionneurs. Le chemin critique le plus long est celui qui est construit par les  $N-1$  additions avec un SCM. Le bloc MCM est la solution développée au niveau bits, la figure V.3 représente l'architecture d'un filtre RIF au forme transposé.

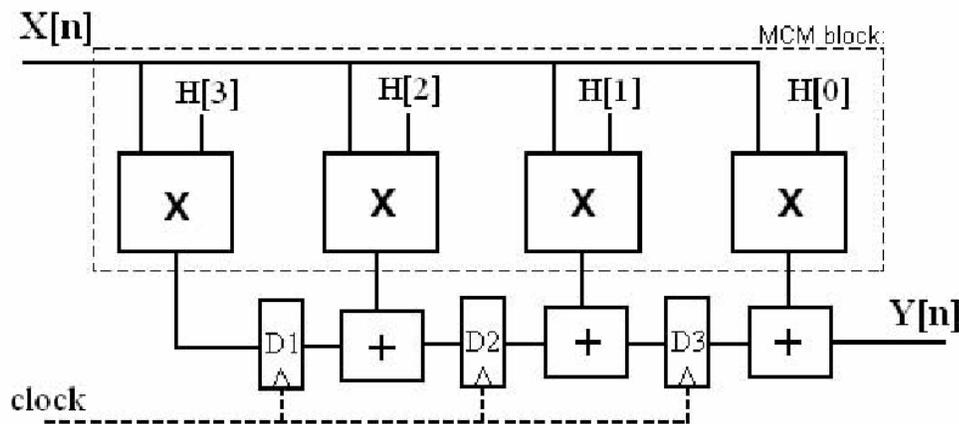


Figure V.3 L'architecture d'un filtre RIF

La génération d'une description matérielle de cette architecture est basée sur une stratégie qui utilise le bloc MCM et des additionneurs pour assurer la sommation des produits (*sum of products*). Pour réduire le chemin critique, augmenter la vitesse des échantillons, et réduire la consommation d'énergie nous utilisons ce qu'on appelle le pipelining.

### V.4.1 Pipelining

Pipelining est une technique de mise en œuvre dans lequel des instructions multiples chevauchent lors de l'exécution pour l'amélioration de la vitesse du chemin critique dans la plus part des systèmes de traitement de signal DSP, microprocesseurs ... etc. Il peut soit augmenter la vitesse d'horloge ou réduire la consommation d'énergie électrique. Le temps d'exécution totale pour chaque instruction individuelle n'est pas modifiée par le pipeline, il n'accélère pas le temps d'exécution des instructions, mais il augmente le nombre d'instructions finies par cycle.

Dans le pipeline toute opération dans le chemin critique est divisée en opérations très petites. L'opération est donc plus rapide avec l'utilisation des registres entre les niveaux, afin d'obtenir un chemin critique plus petit.

**V.5 La conception du programme générateur code VHDL pour les filtres RIF**

Le code VHDL optimisé pour un filtre numérique RIF est basé principalement sur le bloc MCM de notre solution développée précédemment (voire chapitre 4), avec l'addition d'éléments nécessaires à la réalisation du filtre (les registres, l'entrée d'horloge, l'entrée du signal de reset), pour la réinitialisation des registres. Le diagramme ci-après présente les étapes principales de l'algorithme qui génère le filtre RIF optimisé.

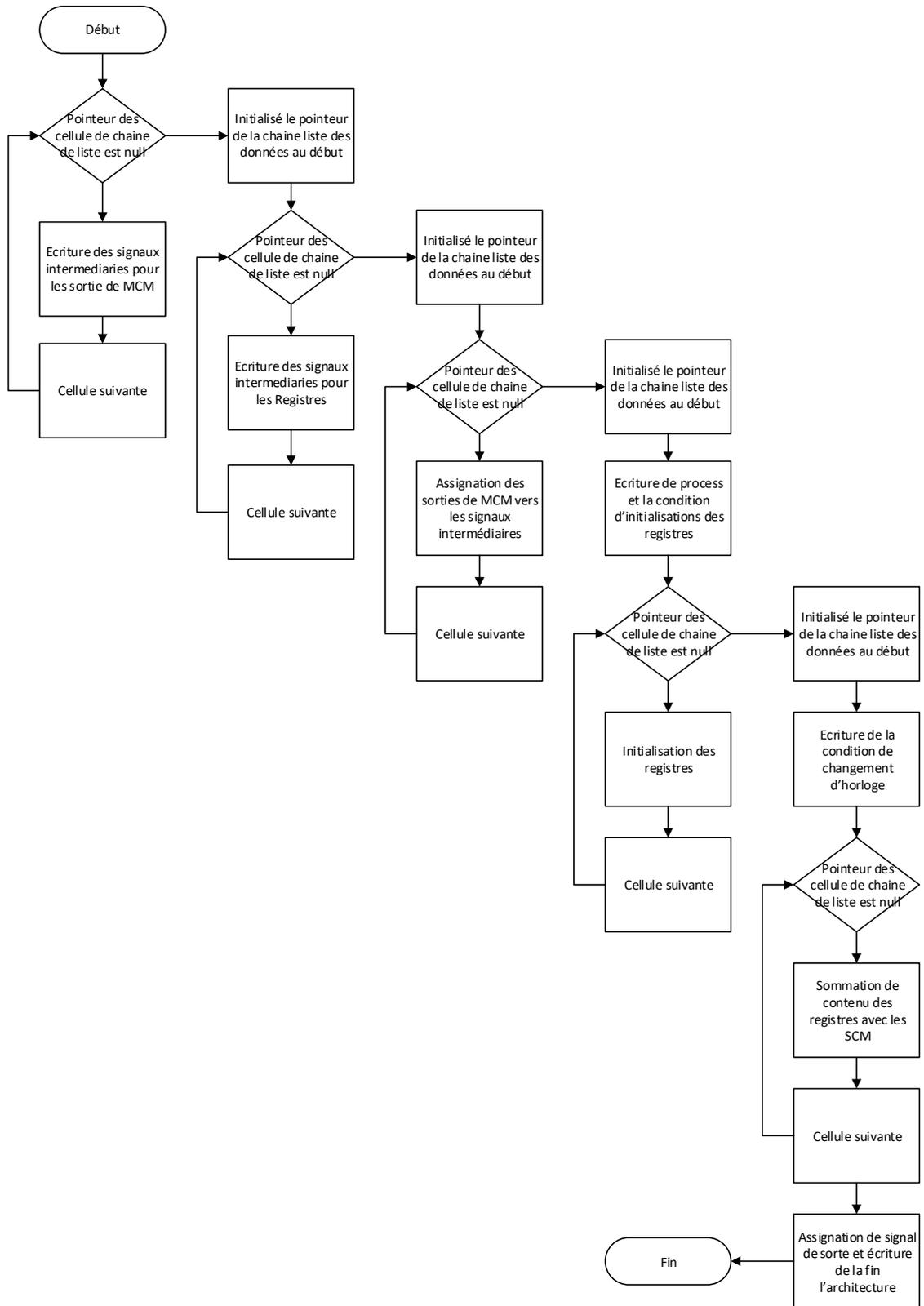


Figure V.4 l'organigramme de generation de description VHDL pour les filtres RIF

## V.5.1 Résultat et discussions

Pour tester l'application, on applique notre programme sur un filtre de 24 coefficients. On extrait le recodage équivalent par le programme RADIX-2<sup>r</sup>, puis nous l'utilisons. Les coefficients du filtre ainsi que son recodage en RADIX-2<sup>r</sup> sont représentés dans la figure V.5.

-710	
327	
505	
582	3 = +1<<1 +1<<0
398	5 = +1<<2 +1<<0
-35	7 = +1<<3 -1<<0
-499	9 = +1<<3 +1<<0
-662	11 = +9<<0 +1<<1
-266	13 = +9<<0 +1<<2
699	35 = +3<<0 +1<<5
1943	133 = +5<<0 +1<<7
2987	199 = +7<<0 +3<<6
3395	291 = +3<<0 +9<<5
2987	327 = +7<<0 +5<<6
1943	331 = +11<<0 +5<<6
699	355 = +3<<0 +11<<5
-266	499 = -13<<0 +1<<9
-662	505 = -7<<0 +1<<9
-499	699 = -5<<0 +11<<6
-35	1943 = -9<<0 -3<<5 +1<<11
398	2987 = +11<<0 -3<<5 +3<<10
582	3395 = +3<<0 +5<<6 +3<<10
505	
327	
-710	



Figure V.5 les coefficients de filtres ainsi que leur recodage en RADIX-2<sup>r</sup>

Après l'exécution de notre programme, il génère une description VHDL optimisé équivalent à partir de recodage précédent (voire la figure V.6).

```

610 SCM_15 <= sig_699 ;
611 SCM_14 <= sig_1943 ;
612 SCM_13 <= sig_2987 ;
613 SCM_12 <= sig_3395 ;
614 SCM_11 <= sig_2987 ;
615 SCM_10 <= sig_1943 ;
616 SCM_9 <= sig_699 ;
617 SCM_8 <= sig_133 & '0';
618 SCM_7 <= sig_331 & '0';
619 SCM_6 <= sig_499 ;
620 SCM_5 <= sig_35 ;
621 SCM_4 <= sig_199 & '0';
622 SCM_3 <= sig_291 & '0';
623 SCM_2 <= sig_505 ;
624 SCM_1 <= sig_327 ;
625 SCM_0 <= sig_355 & '0';
626 filtr : process (clk,reset)
627 begin
628   if reset = '1' then
629
630     MAC_24 <= (OTHERS => '0');
631     MAC_23 <= (OTHERS => '0');
632     MAC_22 <= (OTHERS => '0');
633     MAC_21 <= (OTHERS => '0');
634     MAC_20 <= (OTHERS => '0');
635     MAC_19 <= (OTHERS => '0');
636     MAC_18 <= (OTHERS => '0');
637     MAC_17 <= (OTHERS => '0');
638     MAC_16 <= (OTHERS => '0');
639     MAC_15 <= (OTHERS => '0');
640     MAC_14 <= (OTHERS => '0');
641     MAC_13 <= (OTHERS => '0');

```

Figure V.6 Une partie de la description VHDL de filtre

Pour la simulation de ce filtre sur modelsim, on utilise un fichiers (.do), qui génère des stimuli d'entrée du filtre, le signal de reset, ainsi que le signal d'horloge, comme montrée à la figure V.7.

```

1   force Xb 00000001 0
2   force RESET 0 0,1 10,0 20
3   force CLK 1 0,0 10 -repeat 20
4   force cin 0 0
5   run 300ns

```

Figure V.7 la generation des signaux d'entré pour la simulation

Comme il est apparaît à la figure 5.8, les résultats de simulation du filtre RIF apparaissent après un certain nombre d'oscillation d'horloge, la sortie du filtre sera égale à la somme des produit  $C_i \times X$ , à cause du pipeline du système qui joue le rôle de délai. La partie inconnue dans ce signal (en rouge) est fait parce que les registres n'étaient pas encore remplis par des valeurs, donc on applique une impulsion de reset qui permet la réinitialisation des registres à 0.

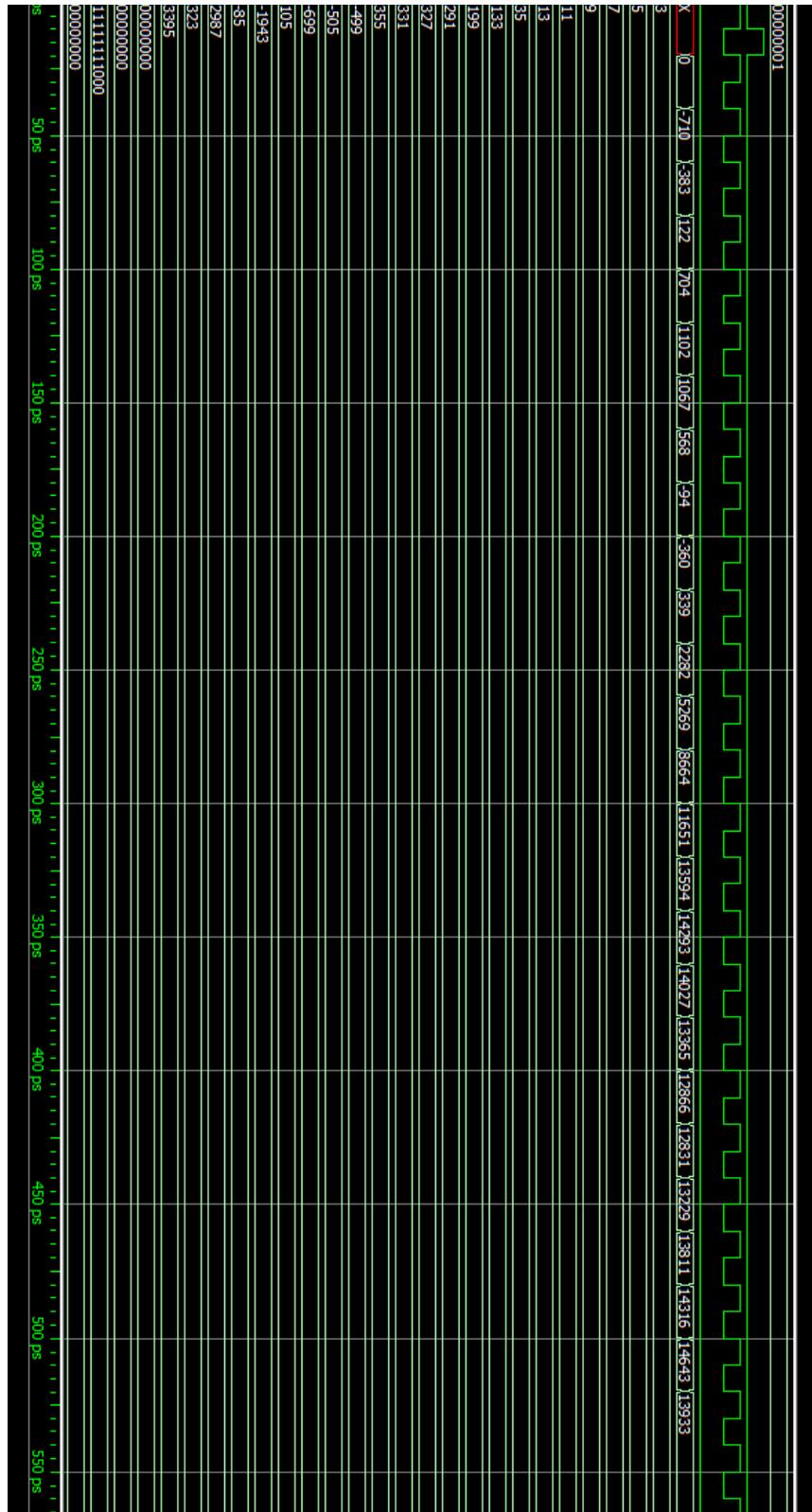


Figure V.8 les résultat de simulation de filtre

## **V.6 Conclusion**

Dans ce chapitre on a appliqué la solution de résolution des problèmes SCM/MCM au niveau d'additionneur de 1 bit, sur les filtres à réponse impulsionnelle finie RIF. Cette application permet de générer une description VHDL très optimisée bit par bit pour tous les filtres RIF. Elle offre une optimisation importante au niveau des ressources consommées dans les blocs de multiplication, avec une architecture interne optimisée basée sur des d'additionneur de 1 bit.

# Conclusion générale

L'objectif de notre projet était le développement d'une nouvelle version de RADIX-2<sup>r</sup> au niveau "Bit" afin de solutionner plus efficacement le problème de la Multiplication par une Constante. Il s'agit de réduire davantage la consommation de puissance, d'améliorer la vitesse, et de réduire la complexité hardware (ressource logiques). Pour cela la solution développée est basée sur l'heuristique RADIX-2<sup>r</sup> avec un algorithme de recodage multi bit .

Nous avons montré que la solution que nous avons développée au niveau bit présente des meilleurs résultats pour l'implémentation hardware de système LTI. La méthodologie adoptée est comme suit:

- Nous avons étudié et comparé les quatre systèmes de recodage : CSD, DBNS, RNS, et Radix-2<sup>r</sup> afin de montrer l'efficacité de ce dernier en terme des performances pour la résolution du problème de la multiplication par une constante
- Nous avons développé une solution au niveau bit de l'heuristique Radix-2<sup>r</sup>. Cette solution est prévisible en termes de nombre maximal d'additionneurs binaires (Full Adder).
- Nous avons montré expérimentalement les performances de notre solution développée au niveau bloc d'additionneur et au niveau bit.
- Nous avons appliqué notre solution sur un filtre FIR afin de générer un filtre FIR optimisé.

Enfin, la nouvelle arithmétique Radix-2<sup>r</sup> peut être avantageusement appliquée à d'autres domaines numériques, tels qu'en traitement numérique du signal (DSP), traitement d'image, télécommunications et le cryptage. Une idée consiste à appliquer la nouvelle heuristique SMC/MCM aux algorithmes de chiffrement RSA pour cibler des longues clés de chiffrement (plus de 4096 bits).

## *Bibliographie*

---

### **Livres:**

- [1] M.D. Ercegovic and T. Lang, "Digital Arithmetic," Morgan Kaufman Publishers, an Imprint of Elsevier Science, ISBN: 1-55860-798-6, San Francisco, USA, © 2004.
- [2] R. Kastner, A. Hosangadi, and F. Fallah, "Arithmetic Optimization Techniques for Hardware and Software Design," Cambridge University Press, ISBN-13 978-0-521-88099-2, © 2010.
- [3] J.M. Muller et al., "Handbook of Floating-Point Arithmetic," Birkhäuser Boston, a part of Springer Science+Business Media, ISBN: 978-0-8176-4704-9, © 2010.
- [4] B. Perhami, "Computer Arithmetic, Algorithms and Hardware Design," Oxford University Press, ISBN: 0-19-512583-5, New-York, USA, © 2000.
- [5] C. Delannoy, "programmeur en langage c cours et exercices corrigés", Paris : ÉDITIONS EYROLLES, 2009, 281p, ISBN 978-2212140101.
- [6] K.C.Chang, "Digital system design with VHDL and synthesis : An integrated approach, Los Alamitos", California: IEEE Computer society, 1999, 551p, ISBN 0-7695-0023-4.
- [7] j. Charles H. Roth, "Digital system design using VHDL, Boston: PWS Publishing" company, 1997, 476p, ISBN 0-534-95099-X.

### **Theses:**

- [8] A. K. Oudjida, «Binary Arithmetic for Finite-Word-Length Linear Controllers: MEMS Applications » Université de Franche-Comté, Besançon, 2015.
- [9] J. Thong, "New Algorithm for Constant Coefficient Multiplication in Custom Hardware," Master Thesis, No 4292, McMaster University, Hamilton, Ontario, Canada, October 2009.
- [10] M. Kamal BOUDJELABA, "contribution à la conception des filtres bidimensionnels non récursifs en utilisant les techniques de l'intelligence artificielle : application au traitement d'images", Thèse de doctorat, Université Ferhat Abbas – Sétif – 1 –, 11 juin 2014.

### **Articles:**

- [11] A. Avizienis, "Signed-digit number representation for fast parallel arithmetic," IRE Trans. on Electronic Computers, vol. EC-10, No. 3, pp. 389–400, September 1961.
- [12] V.S. Dimitrov, G.A. Jullien, and W.C. Miller, "Theory and Applications of the Double-Base Number System," IEEE Trans. on Computers (TC), vol. 48, No. 10, pp. 1098-1106, October 1999.
- [13] H. Sam, and A. Gupta, "A Generalized Multibit Recoding of Two's Complement Binary Numbers and its Proof with Application in Multiplier Implementation," IEEE Trans. on Computers, vol. 39, N° 8, August 1990.

## Bibliographie

---

- [14] J. Thong and N. Nicolici, "An optimal and practical approach to single constant multiplication," *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems*, vol. 30, no. 9, pp. 1373–1386, September 2011.
- [15] Y.E. Kim et al., "Efficient Design of Modified Booth Multipliers for Predetermined Coefficients," *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 27172720, Island of Kos, Greece, May 2006.
- [16] V.S. Dimitrov, L. Imbert, and A. Zakaluzny, "Multiplication by a Constant is Sublinear," *Proceedings of the 18th IEEE Symposium on Computer Arithmetic (ARITH)*, pp. 261-268, June 2007.
- [17] R.I. Hartley, "Subexpression Sharing in Filters Using Canonic Signed Digit Multipliers," *IEEE Trans. on Circuits and Systems II: Analog and Digital Signal Processing*, vol. 43, No. 10, pp. 677-688, October 1996.
- [18] N. Boullis and A. Tisserand, "Some Optimizations of Hardware Multiplication by Constant Matrices," *IEEE Trans. on Computers (TC)*, vol. 54, No. 10, pp. 1271-1282, October 2005.
- [19] R.L. Bernstein, "Multiplication by Integer Constant," *Software– Practice and Experience* 16, 7, pp. 641-652, 1986.
- [20] O. Gustafsson, A.G. Dempster, and L. Wanhammar, "Extended Results for Minimum-Adder Constant Integer Multipliers," *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS)*, vol. 1, pp. I-73 I-76, Scottsdale Arizona, USA, May 2002.
- [21] Y. Voronenko and M. Püschel, "Multiplierless Multiple Constant Multiplication," *ACM Trans. on Algorithms (TALG)*, vol. 3, No. 2, article 11, pp. 1-38, May 2007.
- [22] A. K. Oudjida, N. Chaillet, M. L. Berrandjia «RADIX-2r Arithmetic for Multiplication by a Constant: Further Results and Improvements » *IEEE Trans. on Circuits and Systems II: Express Brief*, vol. 62, pp. 372-376, Avril 2015.
- [23] A. K. Oudjida, N. Chaillet «RADIX-2r Arithmetic for Multiplication by a Constant » *IEEE Trans. on Circuits and Systems II: Express Brief*, vol. 61, pp. 349-353, May 2014.
- [24] A.Liacha, A. K. Oudjida, F. Ferguene, M. Bakiri, M. L. Berrandjia «Design of High-Speed, Low-Power, and Area-Efficient FIR Filters» *Circuits, Devices and Systems*, n° %15, DOI 10.1049/iet-cds.2017.0058 , 2017.
- [25] A. K. Oudjida, A. Liacha, M. Bakiri, N. Chaillet," Multiple Constant Multiplication Algorithm for High-Speed and Low-Power Design" *IEEE Trans on Circuits and Systems II: Express Briefs*, vol. 63, n° %12, pp. 176 - 180, Fevrier 2016.

### Rapport:

## *Bibliographie*

---

- [26] V. Lefèvre, “Multiplication by an Integer Constant,” INRIA Research Report, No. 4192, Lyon, France, May 2001.
- [27] B. Lopez, T. Hilaire and L.S. Didier, “Sum-of-products Evaluation Schemes with Fixed-Point arithmetic, and their application to IIR filter implementation,” Proceedings of the International Conference on Design and Architecture for Signal and Image Processing (DASIP), Karlsruhe, Germany, Oct. 2012.

### **Page web:**

- [28] «The leading operating system for PCs, IoT devices, servers and the cloud \_ Ubuntu,» Canonical Ltd. Ubuntu, 2017. [En ligne]. Available: <https://www.ubuntu.com/>. [Accès le 6 septembre 2017].
- [29] C. f. d. d'Ubuntu, «gedit - Documentation Ubuntu Francophone,» Communauté francophone d'utilisateurs d'Ubuntu, 2017. [En ligne]. Available: <https://doc.ubuntu-fr.org/gedit>. [Accès le 6 septembre 2017].
- [30] C. f. d. d'Ubuntu, «gcc - Documentation Ubuntu Francophone,» Communauté francophone d'utilisateurs d'Ubuntu, 2017. [En ligne]. Available: <https://doc.ubuntu-fr.org/gcc>. [Accès le 6 septembre 2017].
- [31] «Binary Adder and Subtractor,» ELECTRONICS HUB, 29 juin 2015. [En ligne]. Available: <http://www.electronicshub.org/binary-adder-and-subtractor/>. [Accès le 7 septembre 2017].
- [32] CDTA « Multiple Constant Multiplication Algorithm », © 2015. [En ligne]. Available : <http://www.cdta.dz/products/mcm/>. [Accès le 16 septembre 2017].