

9
الجمهورية الجزائرية الديمقراطية الشعبية
REPUBLIQUE ALGERIENNE DEMOCRATIQUE ET POPULAIRE

وزارة التربية الوطنية
MINISTERE DE L'EDUCATION NATIONALE
جامعة قسنطينة
UNIVERSITE DE CONSTANTINE
معهد الإلكترونيك
INSTITUT D'ELECTRONIQUE

FER
2321

THESE DE MAGISTER

Option : Contrôle et Traitement de Signal

Thème

Mise au point d'un système de traitement
numérique du signal en temps réel :
Théorie et Applications

Présentée par :

Abdelhak FERHAT-HAMIDA

Soutenue le 11 Novembre 1992 devant le jury :

MM. L. ABIDA	<i>Président</i>
Y. BOUTERFA	<i>Rapporteur</i>
D. CHIKOUCHE	<i>Rapporteur</i>
A. SAID	<i>Examineur</i>
F. MARIR	<i>Examineur</i>
K. BENMAHAMMED	<i>Examineur</i>

FE R/2321

A mes parents

A mes frères et soeurs

A toute ma famille

REMERCIEMENTS

Je tiens à remercier tous ceux qui ont collaboré de près ou de loin à la réalisation de ce travail. Particulièrement:

Je tiens à exprimer ma profonde reconnaissance à Monsieur Youcef **ROUTERFA**, Professeur à l'Université de Sétif, qui a dirigé ma thèse pour les encouragements et les conseils qu'il m'a prodigués tout au long du travail.

Je remercie Monsieur **L.ZEGADI** pour ses encouragements et son savoir faire monter le moral durant les moments difficiles.

Je remercie Mademoiselle **N.MERGHAD** pour la frappe du manuscrit.

Je remercie Monsieur **L.SELMANI** pour les moyens matériels qu'il a mis à ma disposition durant la réalisation du travail.

Je remercie Monsieur **M.BOUAMAR** pour avoir fourni le système de développement et la documentation de l'ADSP-2100.

Je remercie Monsieur **D.CHIKOUCHE** pour ses conseils.

Je remercie Monsieur **K.BENMAHAMMED** pour la documentation qu'il m'a fournie.

Je remercie mon ami et frère **Fayçal RADJAH** que j'ai trouvé à mes cotés tout au long du travail.

Je remercie vivement les membres de Jury qui ont accepté de juger ce travail:

- Le Rapporteur Monsieur **Y.BOUTERFA**, Professeur à l'Université de Sétif.
- Monsieur **L.ABIDA**, Professeur à l'Université de Batna, qui a accepté de présider le Jury.
- Messieurs les membres **A.SAID** Maître de Conférences à l'Université de Constantine, **F.MARIR** PhD et chargé de cours à l'Université de Constantine, **K.BENMAHAMMED** PhD et maître assistant à l'Université de Blida et **D.CHIKOUCHE** chargé de cours à l'Université de Sétif.

SOMMAIRE

Introduction	1
Chapitre 1: Signaux et systèmes à temps discret	4
1.1. Signaux à temps discret	4
1.1.1. Généralités	4
1.1.2. Norme d'un signal	5
1.1.3. Représentation fréquentielle d'un signal discret	6
1.1.4. Echantillonnage et reconstitution analogique	8
1.2. Systèmes à temps discret	10
1.2.1. Système causal	11
1.2.2. Système linéaire	11
1.2.3. Réponse impulsionnelle d'un système linéaire	11
1.2.4. Système invariant	12
1.2.5. Equation aux différences	12
1.2.6. Système linéaire invariant causal	13
1.2.7. Stabilité	13
1.3. Système de traitement numérique du signal	14
1.4. La transformée en Z	15
1.4.1. Généralité	15
1.4.2. Définition	15
1.4.3. Propriétés de la TZ	15
1.4.4. La transformée en Z inverse	18
1.4.5. Relation avec la transformation de Fourier	19
1.4.6. Fonction de transfert	20
1.5. Signaux discrets aléatoires	21
1.5.1. Introduction	21
1.5.2. Processus aléatoires:	22
1.5.3. Moyennes:	23
1.5.4. Représentation spectrale	26
1.5.5. Réponse des systèmes linéaires aux signaux discrets	28
1.5.6. Bruit blanc	30
Chapitre 2: Transformation de Fourier discrète	32
2.1. Série de Fourier Discrète	32
2.2. Propriétés de la série de Fourier	33
2.2.1. Linéarité	33

2.2.2. Décalage d'une séquence	31
2.2.3. Convolution périodique	33
2.2. Transformation de Fourier Discrète	35
2.3. Propriétés de la TFD :	36
2.3.1. Linéarité	36
2.3.2. Décalage circulaire d'une séquence	36
2.3.3. Propriété de symétrie	37
2.3.4. Convolution circulaire	37
2.5. Transformation de Fourier Discrète pour les signaux à durée illimitée	38
2.6. Convolution linéaire	39
2.7. Transformation de Fourier Rapide FFT	40
2.7.1. Certains aspects de l'algorithme de la FFT	44
2.8. FFT pour les signaux réels	45
2.9. FFT inverse	46
Chapitre 3: Les filtres numériques	48
3.1. Introduction	48
3.2. Filtres IIR	50
3.2.1. Approximations des filtres analogiques	51
3.2.2. Transformation bilinéaire	53
3.2.3. Plan de conception	55
3.3 Filtres FIR	56
3.3.1. Réponse fréquentielle des filtres FIR à phase linéaire	58
3.3.2. Calcul des filtres FIR par développement en série de Fourier	58
3.3.3. Caractéristiques des filtres FIR	59
3.4. Structures des filtres numériques	60
3.4.1. Filtres IIR	60
3.4.2. Filtres FIR	66

Chapitre 4: Effets de la longueur finie des registres en traitement numérique du signal	68
4.1. Introduction	68
4.2. Représentation en virgule fixe	68
4.2.1. Représentation en signe et valeur absolue	69
4.2.2. Représentation en complément à 2	72
4.2.3. Représentation en complément à 1	74
4.3. Effets de la longueur des mots sur les filtres IIR	75
4.3.1. Quantification du signal d'entrée	75
4.3.2. Effets de la quantification du produit des multiplications	77
4.3.3. Les dépassements et mise à l'échelle	81
4.3.4. Les cycles limites	86
4.4. Effets de la longueur des mots sur la FFT	87
4.5. Effets de la longueur des mots sur les filtres FIR	93
4.5.1. Effets de la quantification des résultats des multiplications	93
4.5.2 Prévention contre le débordement	93
Chapitre 5: Implantation des algorithmes de base sur un processeur de signal	94
5.1. Introduction	94
5.2. La FFT	95
5.2.1. La mise conditionnelle en format flottant par bloc	97
5.2.2. La mise inconditionnelle en format flottant par bloc	104
5.2.3. Mise à l'échelle des données à l'entrée	108
5.2.4. Comparaison des trois méthodes	109
5.2.5. La FFT d'une séquence réelle	111
5.2.6. Opération d'inversion de bits	112
5.3. Le filtrage IIR	112
5.3.1. Structure 1D	113
5.3.2. Structure 2D	115
5.3.3. Comparaison des deux structures	115
5.3.4. Temps de calcul	116
5.3.5. Applications	116

5.4. Le filtrage FIR	125
5.4.1. Temps de calcul	125
5.4.2. Applications	125
5.5. Estimation de la fonction d'autocorrélation du bruit de quantification	131
Conclusion	134
Bibliographie	135
ANNEXE 1: Le processeur ASDP-2100 et son système de développement	A.1
ANNEXE 2: Listings des programmes	A.36

INTRODUCTION

Le traitement numérique du signal est l'étude des systèmes et des signaux à l'égard des contraintes imposées par les dispositifs de calcul numérique. Cette discipline a trouvé d'importantes applications dans des domaines tels que le filtrage, les télécommunications, le traitement de la parole, le traitement de l'image et bien d'autres. Et de ce fait, elle est devenue un élément essentiel de la technologie moderne et s'est imposée avec le temps grâce au développement technologique continu.

A l'origine, pour mettre au point les systèmes de traitement, qui étaient analogiques, une réalisation hardware de multiples variantes était nécessaire. Le coût résultant était alors très élevé. Le développement de l'ordinateur a remédié à ce problème en faisant imposer les méthodes de traitement numérique. En effet, par la simulation, il a permis l'étude détaillée des systèmes de traitements analogiques avant leur réalisation matérielle. Ce qui faisait aboutir au système optimal avec un coût considérablement réduit.

Certaines applications ont pu même être exécutées par ordinateurs. Néanmoins, ceux-ci présentaient un inconvénient majeur: les traitements devenant de plus en plus complexes ne pouvaient s'effectuer en temps réel.

Plusieurs recherches ont été lancées pour lever ce problème. Elles ont abouti à la découverte d'algorithmes structurés pouvant être exploités en une conception hardware. Il y a eu, par conséquent,

construction de plusieurs machines câblées offrant de très grandes vitesses d'exécution autour de processeurs généraux.

Après, le progrès de la micro-électronique a donné naissance à des processeurs complexes permettant des calculs très rapides et offrant la flexibilité aux algorithmes et aux formats de données. Puis, il y a eu un impératif de standardisation, ce qui a conduit au développement de processeurs structurés pour satisfaire à une large variété de tâches tels l'INTEL 2920, le TMS320, le DSP56000 et l'ADSP-2100. Ainsi, les processeurs de signal sont devenus des outils de traitement indispensables [1][2].

Le but de ce travail est de mettre au point les algorithmes de base de traitement numérique du signal, à savoir la transformée de Fourier rapide (FFT), le filtrage numérique et la convolution [1][3][4], sur le processeur de signal d'ANALOG DEVICES l'ADSP-2100 pour un traitement en temps réel. Ils seront implantés en virgule fixe simple précision en respectant les deux éléments clés du traitement numérique du signal : la grande vitesse de calcul et la précision numérique adéquate [5].

Ce manuscrit est composé de cinq chapitres et deux annexes:

Le premier chapitre est un rappel sur les signaux et les systèmes discrets, la transformée en Z et les signaux discrets aléatoires. Il contient des notions nécessaires utilisées dans les chapitres suivants.

Le deuxième chapitre introduit la transformée de Fourier discrète. L'algorithme de la FFT à entrelacement temporel et son application pour les signaux réels sont expliqués.

Le troisième chapitre traite les deux classes de filtres numériques à savoir les filtres IIR et les filtres FIR. Pour chaque classe, une méthode d'approximation et les structures les plus utilisées sont données.

Dans le quatrième chapitre, on présente l'effet de la longueur finie des registres sur les algorithmes développés aux deux chapitres précédents.

Le dernier chapitre est consacré à l'étude de l'implantation des algorithmes sur le processeur de signal ADSP-2100. Les résultats obtenus et quelques applications y sont présentés.

Enfin, dans les annexes, on présente le processeur et son système de développement ainsi que les listings des programmes développés en langage assembleur.

1.1. Signaux à temps discret:1.1.1 Généralités :

Un signal est une grandeur physique qui évolue avec le temps et qu'on peut représenter mathématiquement par une fonction d'une ou de plusieurs variables indépendantes. Si la représentation mathématique est déterminée d'une manière unique par une loi qu'on peut énoncer, le signal est dit déterministe (figure 1.1.a), sinon il sera régi par les lois de probabilité et le signal est dit aléatoire (figure 1.1.b) [6-8][11].

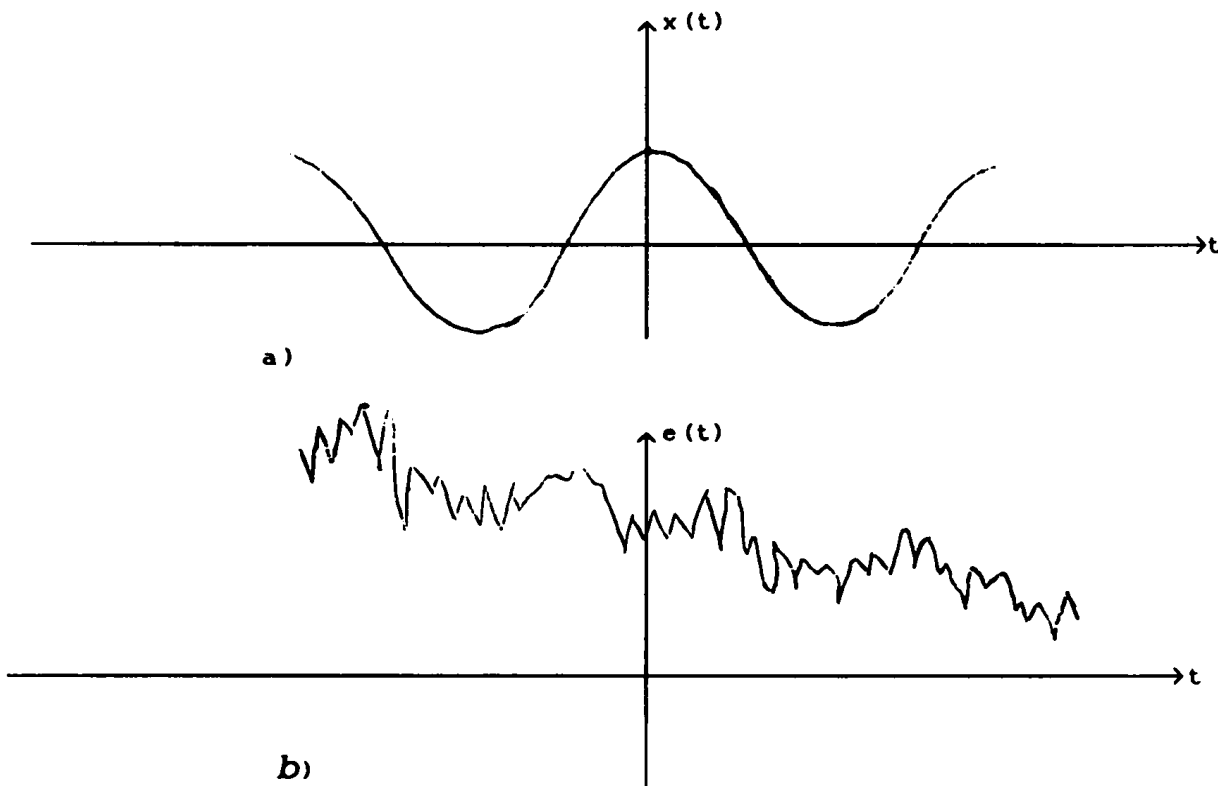


Figure 1.1 Types de signaux

- a) Déterministe- signal cosinus
- b) Aléatoire

La variable indépendante peut être continue ou discrète.

- Les signaux à temps continu, ou signaux analogiques, sont définis dans un continuum de temps et sont ainsi représentés par les fonctions à variables continues (figure 1.2.a) [7].

- Les signaux à temps discret sont ceux pour lesquels la variable indépendante prend uniquement des valeurs discrètes, généralement espacées d'un intervalle de temps constant T appelé période d'échantillonnage [6][7][11]. Les signaux discrets sont représentés par des séquences de nombres (figure 1.2.b) [6]:

$\{ x(nT) \}, n = \{ \dots, -3, -2, -1, 0, 1, 2, 3, \dots \}$ (1.1)
 $x(nT)$ est le $n^{\text{ème}}$ membre de la séquence. Par simplicité, on le dénote $x(n)$, $T = 1$

- L'amplitude du signal peut être aussi continue ou discrète. Les signaux digitaux (ou signaux numériques) sont ceux pour lesquels le temps et l'amplitude sont discrets (figure 1.2.c) [7].

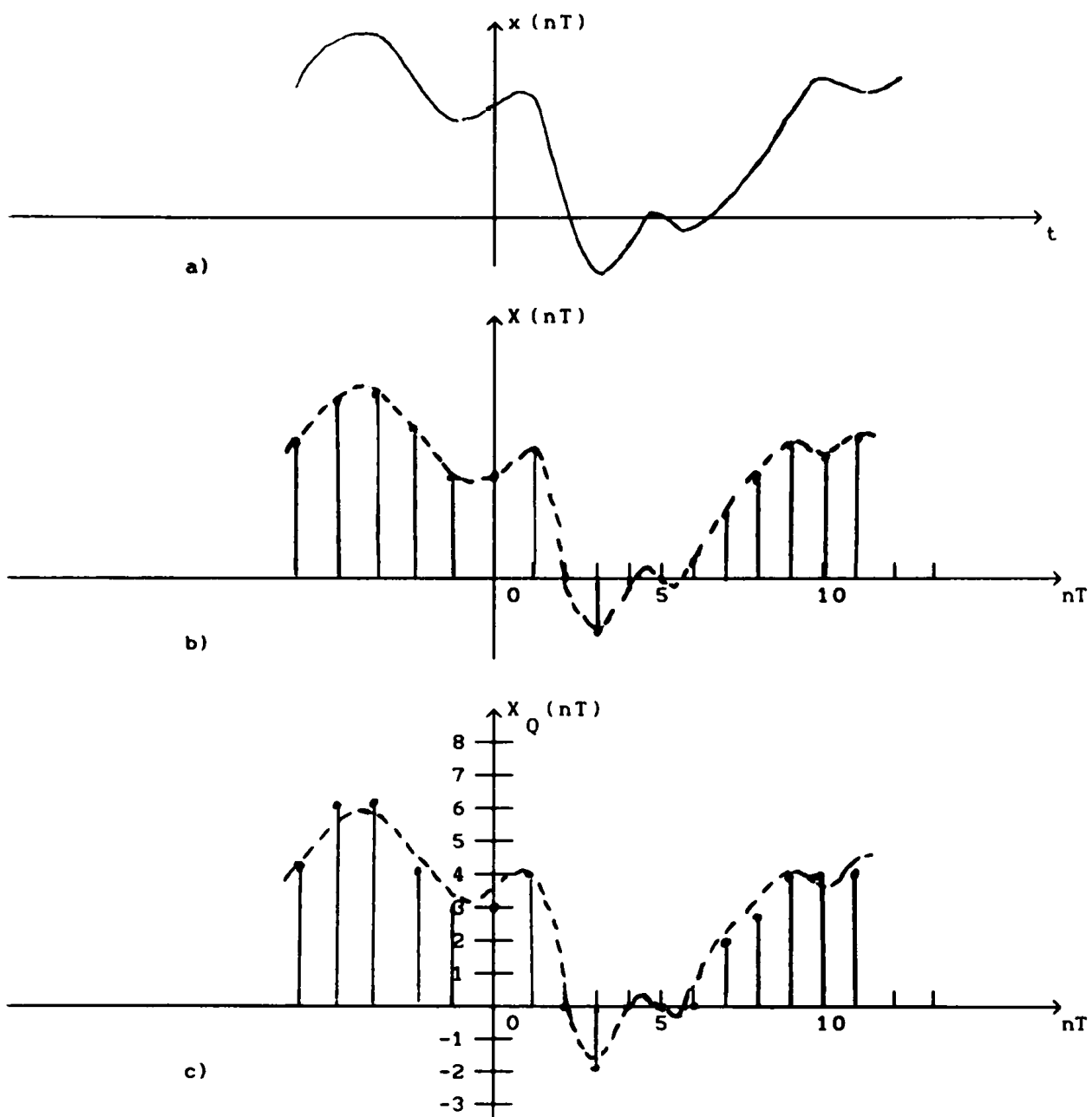


Figure 1.2 Types de signaux déterministes:
 a) analogique
 b) discret
 c) digital (numérique)

Dans tout ce qui suit, on considère que les signaux sont discrets complexes et que la période d'échantillonnage est égale à l'unité. Les exceptions à ceci seront précisées quand il le faut.

1.1.2 Norme d'un signal :

La norme L_p d'un signal est définie par [6]:

$$L_p = || X(n) ||_p = \left| \sum_{n=-\infty}^{+\infty} |x(n)|^p \right|^{1/p} \quad (1.2)$$

où p est un entier positif. La norme d'un signal fournit une mesure globale de sa dimension.

La norme d'un signal satisfait aux axiomes suivants :

- 1- $|| x(n) || \geq 0$ et $|| x(n) || = 0$ si et seulement si $x(n) = 0$ pour tout n.
- 2- $|| a \cdot x(n) || \leq |a| \cdot || x(n) ||$, pour tout scalaire a.
- 3- $|| x(n) + y(n) || \leq || x(n) || + || y(n) ||$

Trois normes portent un intérêt particulier en traitement du signal: L_1 , L_2 , L_∞

- La norme L_1 est égale à la somme des amplitudes de chaque échantillon du signal:

$$L_1 = || x(n) ||_1 = \sum_{n=-\infty}^{+\infty} |x(n)| \quad (1.3)$$

Cette mesure est utilisée dans la détermination de la stabilité des systèmes discrets linéaires.

- La norme L_2 est une mesure de l'énergie du signal.

$$L_2 = || x(n) ||_2 = \left(\sum_{n=-\infty}^{+\infty} |x(n)|^2 \right)^{1/2} \quad (1.4)$$

Elle est utilisée dans l'analyse des signaux et des systèmes.

- La norme L_∞ donne le maximum de l'amplitude du signal :

$$L_\infty = || x(n) ||_\infty = \max |x(n)|, \text{ pour tout } n. \quad (1-5)$$

Cette norme fournit une limite pour déterminer les exigences de la dynamique des systèmes.

1.1.3 Représentation fréquentielle d'un signal discret :

La transformée de Fourier d'un signal analogique $x_a(t)$ est [7]:

$$X_a(\omega) = \int_{-\infty}^{+\infty} x_a(t) \cdot e^{-j\omega t} dt \quad (1.6)$$

et la transformée inverse de $X_a(\omega)$ est [7]:

$$x_a(t) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} X_a(\omega) \cdot e^{j\omega t} d\omega \quad (1.7)$$

où ω dénote la pulsation et t dénote le temps.

La transformée de Fourier du signal discret $\{x(n)\}$ est donnée par [7]:

$$X(\omega) = \sum_{n=-\infty}^{+\infty} x(n) \cdot e^{-j\omega n} \quad (1.8)$$

De cette équation, on déduit que $X(\omega)$ est périodique de période 2π , puisqu'on a :

$$e^{j(\omega+2\pi)n} = e^{j\omega n}$$

L'équation (1.8) exprime $X(\omega)$ sous la forme d'une série de Fourier où les échantillons $x(n)$ correspondent aux coefficients de la série. De ce fait, on peut évaluer les échantillons par la relation donnant les coefficients de Fourier d'une fonction périodique:

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{+\pi} X(\omega) \cdot e^{j\omega n} d\omega \quad (1.9)$$

c'est la transformée de Fourier inverse de $X(\omega)$.

Existence de la transformée de Fourier:

On admet que la transformée de Fourier définie par (1.8) existe pour les signaux à énergie finie, ou de carré sommable, i.e tous les signaux qui vérifient la relation [8]:

$$\sum_{n=-\infty}^{+\infty} |x(n)|^2 < \infty$$

Notations fréquentielles:

Certains auteurs [8] préfèrent travailler avec la variable f représentant la fréquence à la place de la variable ω représentant la pulsation. On a $\omega = 2\pi f$. Les équations (1.8) et (1.9) deviennent respectivement:

$$X(f) = \sum_{n=-\infty}^{+\infty} x(n) \cdot e^{-j2\pi f n}$$

et

$$x(n) = \int_{-1/2}^{+1/2} X(f) \cdot e^{j2\pi f n} df$$

Dans ce qui suit, on utilisera indifféremment les deux notations.

La transformée de Fourier et le produit de convolution:

Soient deux signaux $x(n)$ et $y(n)$ ayant pour transformées respectives $X(\omega)$ et $Y(\omega)$, et soit $z(n)$ leur produit de convolution défini par [7]:

$$z(n) = x(n) * y(n) = \sum_{k=-\infty}^{+\infty} x(k) \cdot y(n-k) \quad (1.10)$$

La transformée de Fourier de $z(n)$ est donnée par le produit simple:

$$Z(\omega) = X(\omega) \cdot Y(\omega) \quad (1.11)$$

1.1.4 Echantillonnage et reconstitution analogique :

Souvent, les signaux discrets sont obtenus par un échantillonnage périodique d'un signal analogique. Les échantillons sont prélevés avec une période T appelée période d'échantillonnage. Dans le cas où on désire reconstituer le signal analogique à partir de ses échantillons, on doit poser une contrainte sur T , cette contrainte découle du théorème d'échantillonnage.

1.1.4.1 Théorème d'échantillonnage : Théorème de Shannon

Théorème: Un signal analogique $x_a(t)$ ayant une largeur de bande finie limitée à Ω_c ne peut être reconstitué exactement à partir de ses échantillons $x_a(kT)$ que si ceux-ci ont été prélevés avec une période T inférieure ou égale à (π/Ω_c) [8].

1.1.4.2 Effet de l'échantillonnage :

L'échantillonnage idéal est obtenu en multipliant le signal analogique $x_a(t)$ par une suite périodique d'impulsions de Dirac de période T [8]:

$$x_e(t) = x_a(t) \cdot \delta_T(t)$$

avec

$$\delta_T(t) = \sum_{n=-\infty}^{+\infty} \delta(t-nT)$$

Dans le domaine fréquentiel, la transformée de Fourier $X_e(\omega)$ du signal échantillonné est donnée par:

$$X_e(\omega) = X_a(\omega) * E(\omega) = \frac{1}{T} \sum_{n=-\infty}^{+\infty} X_a\left(\omega - \frac{2\pi n}{T}\right) \quad (1.12)$$

où $E(\omega)$ est la transformée de Fourier de $\delta_T(t)$:

$$E(\omega) = \frac{2\pi}{T} \sum_{n=-\infty}^{+\infty} \delta\left(\omega - \frac{n2\pi}{T}\right)$$

La formule (1.12) montre que le spectre $X_e(\omega)$ est obtenu par la répartition périodique, de période $2\pi/T$, du spectre $X_a(\omega)$. Si le théorème de Shannon n'est pas respecté, lors de la répartition de $X_a(\omega)$, on aura un recouvrement et la reconstitution du signal original n'est plus possible. La figure (1.3) montre deux cas d'échantillonnage, l'un permet la reconstitution et l'autre ne le permet pas.

1.1.4.3 Reconstitution :

Si le théorème d'échantillonnage est satisfait, le signal $x_a(t)$ peut être reconstitué avec un filtre passe-bas idéal. On a ainsi [8]:

$$X_a(\omega) = X_e(\omega) \cdot H(\omega) \quad (1.13)$$

où $H(\omega)$ est la réponse du filtre. Dans le domaine temporel, on a :

$$x_a(t) = x_e(t) * h(t) \quad (1.14)$$

$h(t)$ étant la réponse impulsionnelle du filtre. Elle est donnée par :

$$h(t) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} H(\omega) e^{j\omega t} d\omega = \frac{T}{2\pi} \int_{-\omega_c}^{+\omega_c} e^{j\omega t} d\omega = \frac{\sin(\pi t/T)}{(\pi t/T)} \quad (1.15)$$

ω_c étant la fréquence de coupure du filtre.

En remplaçant (1.15) dans (1.14), on obtient une formule d'interpolation pour la reconstitution du signal analogique $x_a(t)$.

$$x_a(t) = \sum_{k=-\infty}^{+\infty} x_e(k.T) \frac{\sin [(\pi/T) \cdot (t-kT)]}{(\pi/T) \cdot (t-kT)} \quad (1.16)$$

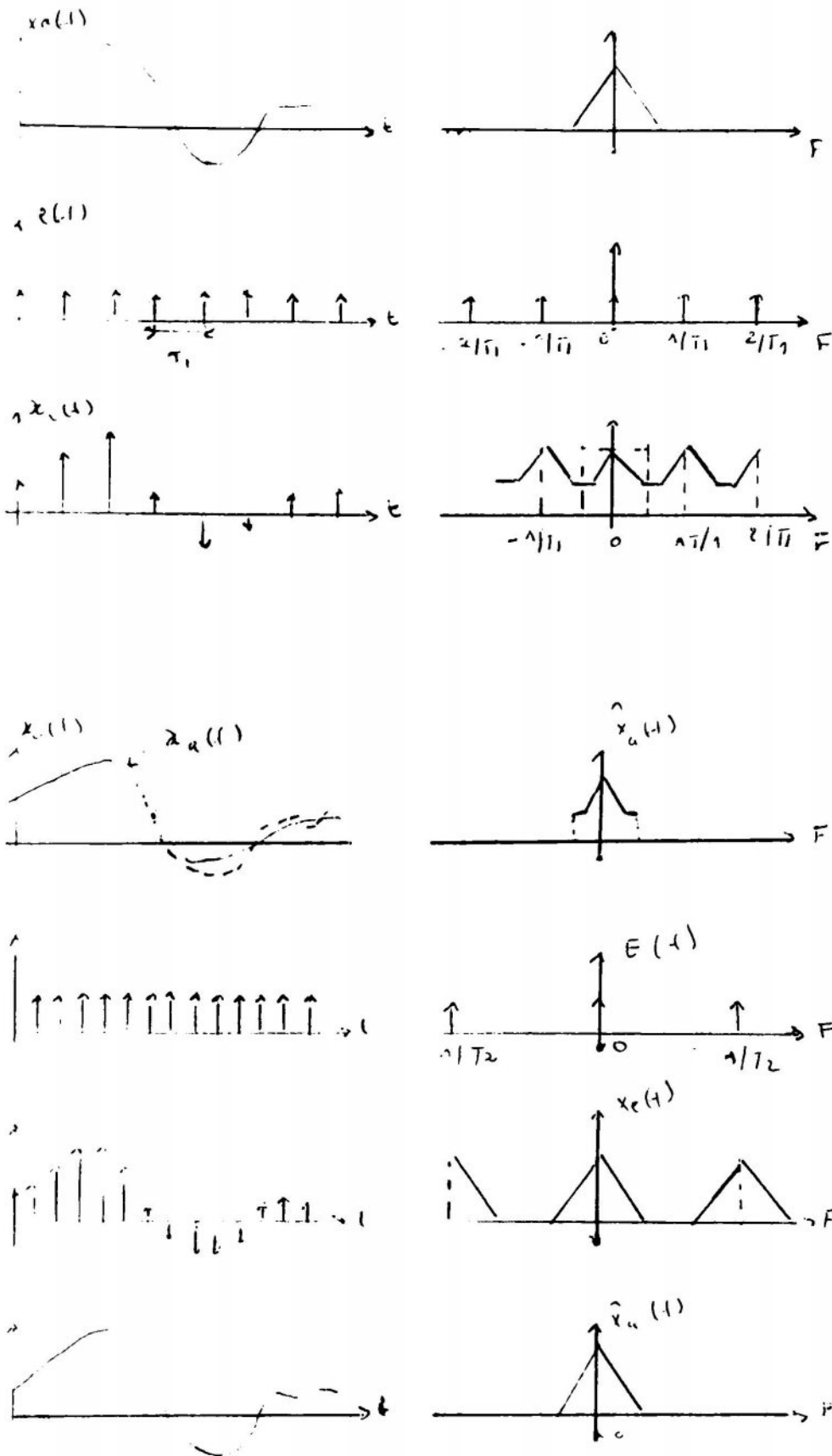


Figure 1.3: Echantillonnage et reconstitution du signal.

1.2 Systèmes à temps discret :

Un système est défini mathématiquement par une transformation unique, ou un opérateur T , qui agit sur une séquence d'entrée $x(n)$ pour produire une séquence de sortie $y(n)$. Ceci est dénoté

par [7]:

$$y(n) = T [x(n)] \quad (1.17)$$

Le signal d'entrée est appelé signal d'excitation et le signal de sortie est appelé la réponse du système à l'excitation.

Les différentes classes de systèmes sont obtenues en posant des contraintes sur l'opérateur T.

1.2.1 Système causal :

Un système causal est caractérisé par le fait que sa réponse ne précède jamais son excitation, c'est à dire, si on a :

$$x(n) = 0 \quad \text{pour} \quad n < n_0 \quad (1.18)$$

alors, on doit avoir :

$$y(n) = 0 \quad \text{pour} \quad n < n_0 \quad (1.19)$$

$y(n)$ étant la réponse à l'excitation $x(n)$.

Pour ces systèmes, l'opérateur T est contraint à ne pas dépendre des valeurs futures de l'excitation.

1.2.2 Système linéaire :

Pour les systèmes linéaires, on pose sur l'opérateur T la contrainte du principe de superposition. On doit avoir :

$$\begin{aligned} T [ax_1(n) + bx_2(n)] &= a T[x_1(n)] + b T[x_2(n)] \\ &= a y_1(n) + b y_2(n) \end{aligned} \quad (1.20)$$

où a et b sont des constantes et $y_1(n)$ et $y_2(n)$ sont les réponses respectives aux excitations $x_1(n)$ et $x_2(n)$.

1.2.3. Réponse impulsionnelle d'un système linéaire :

Une séquence arbitraire peut être exprimée par la somme pondérée d'impulsions unité décalées [7]:

$$x(n) = \sum_{k=-\infty}^{+\infty} x(k) \cdot \delta(n-k) \quad (1.21)$$

Soit $h_k(n)$ la réponse du système à l'impulsion $\delta(n-k)$. Alors des équations (1.20) et (1.21), on obtient :

$$\sum_{n=0}^N a_n \cdot y(k-n) = \sum_{m=0}^M b_m \cdot x(k-m) \quad (1.27)$$

Cette relation permet d'obtenir le kⁱème échantillon de la réponse :

$$y(k) = \sum_{m=0}^M \frac{b_m}{a_0} x(k-m) - \sum_{n=1}^N \frac{a_n}{a_0} y(k-n) \quad (1.28)$$

1.2.6 Système linéaire invariant causal :

Un système linéaire invariant est causal si sa réponse impulsionnelle s'annule pour les arguments négatifs de k [8]:

$$h(k) = 0 \quad \text{pour } k < 0 \quad (1.29)$$

Le produit de convolution donné par (1.25) devient :

$$y(n) = \sum_{k=0}^{+\infty} h(k) \cdot x(n-k) \quad (1.30)$$

D'une manière générale, un signal x(n) est causal si on a:

$$x(k) = 0 \quad \text{pour } k < 0 \quad (1.31)$$

En pratique, seuls les systèmes et les signaux causaux sont pratiquement réalisables.

1.2.7 Stabilité :

Un système est stable si à une excitation bornée, il produit une réponse bornée.

Pour les systèmes invariants :

$$\text{Si on a} \quad |x(n)| \leq M$$

on doit avoir :

$$|y(n)| = \left| \sum_{k=-\infty}^{+\infty} h(k) \cdot x(n-k) \right| \leq M \sum_{k=-\infty}^{+\infty} |h(k)| < \infty$$

Cette relation est vraie si on a :

$$\sum_{k=-\infty}^{+\infty} |h(k)| < \infty \quad (1.32)$$

Les systèmes linéaires invariants sont stables si la relation (1.32) est vérifiée.

1.3. Système de traitement numérique du signal.

Un système de traitement numérique du signal est illustré à la figure (1.4).

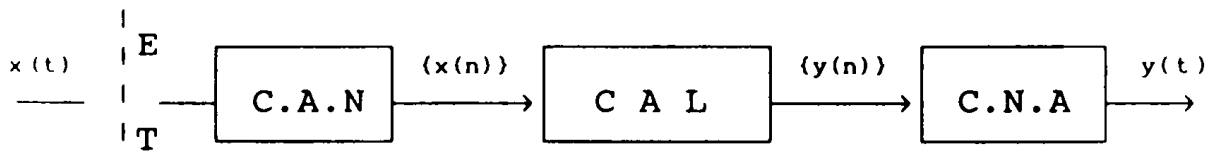


Figure 1.4. Systeme de traitement numerique du signal

Le signal analogique est prélevé toutes les T secondes par un échantillonneur (E). Il est ensuite converti sous forme numérique par un convertisseur analogique-numérique (C.A.N). Le traitement proprement dit est effectué par un ordinateur (C A L) qui, en principe, communique ses résultats toutes les T secondes à un convertisseur numérique-analogique (C.N.A) qui reconstitue le signal à temps continu.

Lorsque le traitement opéré par le ordinateur se fait au même rythme que l'échantillonnage, on dit qu'il s'agit d'un *traitement en temps reel*. Sinon, on est en présence d'un *traitement en temps differe*.

Le traitement en temps réel exige le plus souvent des dispositifs numériques spécialisés ayant une puissance de calcul suffisante pour élaborer les échantillons de sortie à la cadence requise [11]. La partie ordinateur de tels système est souvent un processeur de signal.

Une grande variété de processeurs a vu le jour grâce au progrès continu de la micro-électronique d'une part, et aux exigences de plus en plus sévères des applications d'autre part.

L'ADSP-2100 de la compagnie ANALOG DEVICES est un de ces processeurs. Il est conçu et optimisé pour les traitements en temps réel nécessitant des vitesses de calcul très grandes. La description complète de ce processeur et de son système de développement est donné dans l'annexe 1.

1.4. La transformée en Z.

1.4.1 Généralités:

Dans la théorie des systèmes continus, la transformée de Laplace est considérée comme la généralisation de la transformée de Fourier. Elle joue un rôle important dans l'analyse et la synthèse de ces systèmes. Le besoin est de créer un outil puissant qui jouerait pour les systèmes discrets le même rôle joué par la transformée de Laplace pour les systèmes continus. Ce besoin est comblé par la transformée en Z.

1.4.2. Définition:

La transformée en Z, $X(z)$, d'une séquence $x(n)$ est définie par [6]:

$$\text{TZ } [x(n)] = X(z) = \sum_{n=-\infty}^{+\infty} x(n) \cdot z^{-n} \quad (1.33)$$

où z est une variable complexe et TZ est l'opérateur définissant la transformée.

Cette relation est connue sous le nom de transformée en Z bilatérale. Dans l'analyse et la synthèse des systèmes causaux, on utilise la transformée en Z unilatérale définie par [6]:

$$X(z) = \sum_{n=0}^{+\infty} x(n) \cdot z^{-n} \quad (1.34)$$

Par la suite, on étudiera la transformée en Z bilatérale. La transformée unilatérale en sera un cas particulier.

1.4.3. Propriétés de la TZ:

1.4.3.1 Région de convergence:

$X(z)$ est donnée sous la forme d'une série de puissances et elle pose ainsi le problème de la convergence. On définit la région de convergence comme l'ensemble des valeurs de z pour lesquelles la série (1.33) converge.

Pour trouver la région de convergence, on peut appliquer le critère de Cauchy. Ce critère affirme qu'une série du type:

$$\sum_{k=0}^{+\infty} U_k = U_0 + U_1 + U_2 + \dots \quad (1.35)$$

converge si la condition suivante est satisfaite :

$$\lim_{k \rightarrow \infty} |U_k|^{1/k} < 1 \quad (1.36)$$

On trouve que la série (1.33) converge dans un anneau du plan complexe des z donné par [8]:

$$0 \leq R_{x-} < |z| < R_{x+} \leq +\infty \quad (1.37)$$

Ceci est illustré sur la figure (1.5).

Il est évident que si $R_{x-} > R_{x+}$, $X(z)$ ne converge pas.

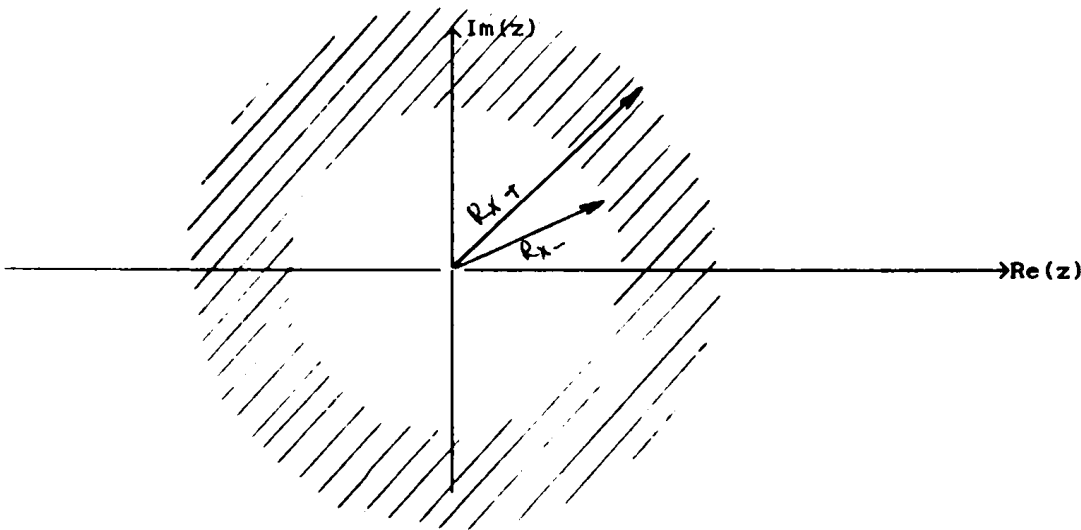


Figure 1.5 : Region de convergence d'une sequence bilaterale

Propriétés :

De l'étude de la convergence, on peut tirer les propriétés suivantes [8]:

- La transformée en Z d'un signal causal converge à l'extérieur d'un cercle de rayon R_{x-} (Fig 1.6.a).

- Similairement, la transformée en Z d'un signal qui est nul pour $n \geq 0$, qu'on appelle signal anti-causal converge à l'intérieur d'un cercle de rayon R_{x+} (Fig 1.6.b).

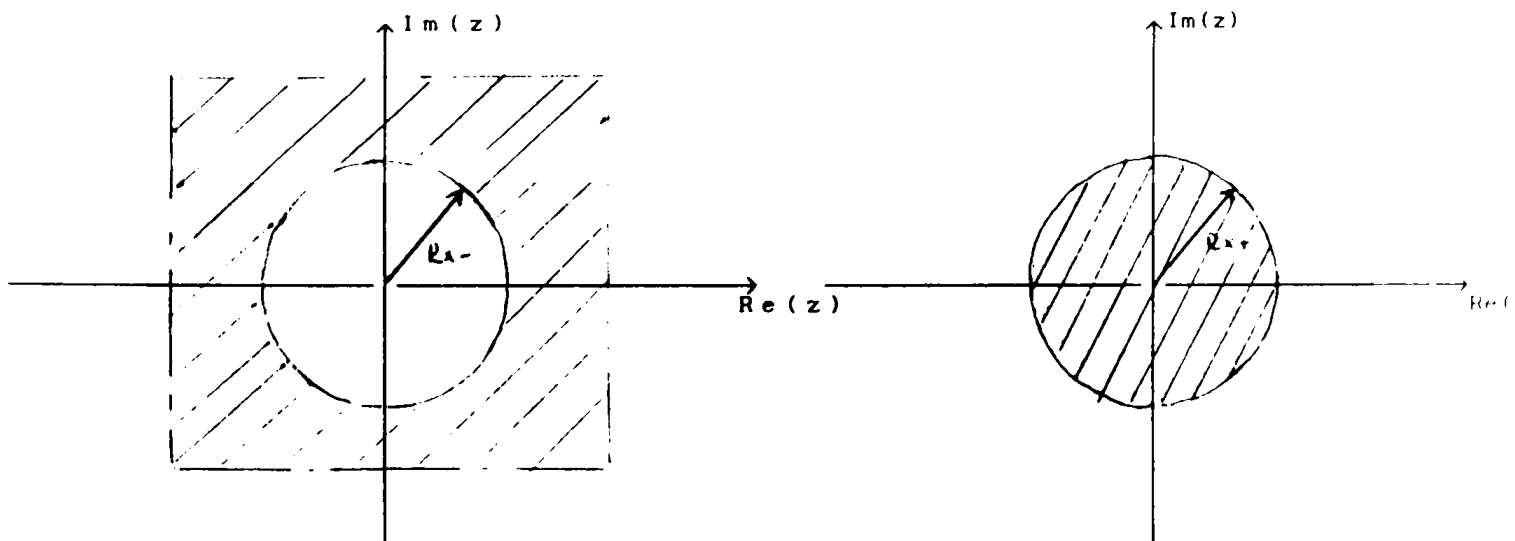


Figure 1.6 : Region de convergence
 a) D'un signal causal
 b) D'un signal anti-causal

Pour un signal à durée finie, la transformée en Z est [8]:

$$X(z) = \sum_{n=n_1}^{n_2} x(n) \cdot z^{-n} \quad (1.38)$$

où n_1 et n_2 sont des entiers. Cette série converge partout sauf peut être pour $z = 0$ et $z \rightarrow \infty$. Pour ces deux valeurs, on peut distinguer trois cas :

- Si n_1 et n_2 sont tous les deux positifs, la série (1.38) ne converge pas pour $z = 0$ car pour $n > 0$, le terme z^{-n} diverge.
- Si n_1 est négatif et n_2 est positif, la série (1.38) diverge pour $z = 0$ et $|z| \rightarrow +\infty$
- Si n_1 et n_2 sont tous les deux négatifs, la série diverge pour $|z| \rightarrow +\infty$

1.4.3.2 Linéarité:

Soient $X(z)$ et $Y(z)$ les transformées en Z des séquences $x(n)$ et $y(n)$. Alors, on a [7]:

$$TZ [ax(n) + by(n)] = aX(z) + bY(z) \quad (1.39)$$

La région de convergence de la transformée de la somme est l'intersection des régions de convergences de la transformée des deux séquences.

1.4.3.3. Décalage d'un signal:

Soient $X(z)$ la transformée en Z de la séquence $x(n)$. La transformée du signal $x(n+n_0)$, où n_0 est un entier, est donnée par [7]:

$$\text{TZ } [x(n+n_0)] = z^{n_0} X(z) \quad (1.40)$$

Pour un signal causal, on distingue deux cas :

- Si n_0 est négatif, la transformée en Z de la version décalée du signal est donnée par l'expression (1.40).
- Si n_0 est positif, cette transformée devient [9]:

$$\text{TZ } [x(n+n_0)] = z^{n_0} \left[X(z) - \sum_{m=0}^{n_0-1} x(m) \cdot z^{-m} \right] \quad (1.41)$$

1.4.3.4. Convolution de séquence :

Le produit de convolution de deux signaux $x(k)$ et $y(k)$ est donné par [7]:

$$w(n) = \sum_{k=-\infty}^{+\infty} x(k) \cdot y(n-k) \quad (1.42)$$

La transformée en Z de ce signal est donnée par :

$$W(z) = X(z) \cdot Y(z) \quad (1.43)$$

La région de convergence de $W(z)$ peut être plus large que l'intersection des régions de convergence de $X(z)$ et $Y(z)$ si les zéros de l'une des transformées compensent les pôles de l'autre.

1.4.4. La transformée en Z inverse :

La relation de la transformée en Z inverse peut être obtenue en utilisant le théorème de Cauchy sur l'intégration le long d'un contour dans le plan complexe. Ce théorème donne [7]:

$$\frac{1}{2\pi j} \oint_C z^{k-1} dz = \begin{cases} 1, & k = 0 \\ 0, & k \neq 0 \end{cases} \quad (1.44)$$

où C est un contour fermé entourant l'origine du plan des z .

En multipliant les deux membres de l'équation (1.33) par $z^{k-1}/2\pi j$ et en intégrant le long d'un contour entourant l'origine et contenu dans la région de convergence, on obtient :

$$\frac{1}{2\pi j} \oint_C X(z) \cdot z^{k-1} dz = \sum_{n=-\infty}^{+\infty} x(n) \cdot \frac{1}{2\pi j} \oint_C z^{-n+k-1} dz \quad (1.45)$$

En utilisant la relation (1.44), on obtient :

$$\frac{1}{2\pi j} \oint_C X(z) \cdot z^{n-1} dz = x(n)$$

Soit la transformée en Z donnée par l'intégrale [7]:

$$x(n) = \frac{1}{2\pi j} \oint_C X(z) \cdot z^{n-1} dz \quad (1.46)$$

Pour la transformée en Z rationnelle, l'intégrale (1.46) est évaluée en utilisant le théorème des résidus [7]:

$$\begin{aligned} x(n) &= \frac{1}{2\pi j} \oint_C X(z) \cdot z^{n-1} dz \\ &= \sum [\text{résidus de } X(z) \cdot z^{n-1} \text{ aux pôles à l'intérieur} \\ &\quad \text{du contour C}] \end{aligned} \quad (1.47)$$

Le résidu à un pôle z_0 d'ordre q est donné par [8]:

$$\begin{aligned} \text{Résidu } [X(z) \cdot z^{n-1}]_{z=z_0} &= \\ \lim_{z \rightarrow z_0} \frac{1}{(q-1)!} \frac{d^{q-1}}{dz^{q-1}} [X(z) \cdot z^{n-1} (z-z_0)^q] & \end{aligned} \quad (1.48)$$

1.4.5. Relation avec la transformation de Fourier :

Dans l'expression (1.33), on peut représenter z à l'aide des coordonnées polaires dans le plan complexe [8]:

$$z = r \exp(j\omega) \quad (1.49)$$

En substituant cette relation dans (1.33), on obtient :

$$X(re^{j\omega}) = \sum_{n=-\infty}^{+\infty} x(n) \cdot r^{-n} \cdot e^{-j\omega n} \quad (1.50)$$

En comparant cette relation avec la définition de la transformée de Fourier (1.12), on remarque que la transformée en Z

s'identifie à la transformée de Fourier pour $|z| = 1$, c'est à dire pour $r = 1$.

$$X(z) \Big|_{|z|=1} = X(\omega) \quad (1.51)$$

1.4.6. Fonction de transfert :

On sait que le signal de sortie d'un système linéaire invariant ayant une réponse impulsionnelle $h(n)$ est donné par le produit de convolution :

$$y(n) = \sum_{k=-\infty}^{+\infty} h(k) \cdot x(n-k) \quad (1.52)$$

et d'après la relation (1.43), on obtient :

$$Y(z) = H(z) \cdot X(z) \quad (1.53)$$

La fonction $H(z)$ est appelée fonction de transfert. Si elle est évaluée pour $|z| = 1$, on obtient la réponse fréquentielle $G(\omega)$ du système.

1.4.6.1 Cas du système causal :

Pour un système causal, la fonction de transfert est donnée par :

$$H(z) = \sum_{n=0}^{+\infty} h(n) \cdot z^{-n} \quad (1.54)$$

Conformément au paragraphe (1.4.3.1), $H(z)$ converge à l'extérieur d'un cercle de rayon R_h donné par :

$$R_h = \lim_{n \rightarrow \infty} |h(n)|^{1/n} \quad (1.55)$$

Comme une transformée en Z ne converge jamais à un pôle, tous les pôles de $H(z)$ doivent être à l'intérieur de ce cercle.

1.4.6.2. Cas du système stable :

Si le système est stable, on doit avoir :

$$\sum_{n=-\infty}^{+\infty} |h(n)| < \infty \quad (1.56)$$

La fonction de transfert $H(z)$ est de la forme (1.33), et elle converge dans un anneau du plan des z . Or, comme la condition (1.56) est satisfaite, $H(z)$ converge pour $|z| = 1$. Par conséquent, l'anneau de convergence doit nécessairement contenir le cercle unité [8].

1.4.6.3 Cas d'un système causal et stable :

Si le système linéaire invariant est causal et stable, le cercle unité doit être contenu dans la région de convergence d'un système causal. Par conséquent, pour un système causal et stable, la fonction de transfert converge à l'extérieur d'un cercle dont le rayon est dans tous les cas inférieur à l'unité.

Ainsi, les pôles de la fonction de transfert d'un système linéaire invariant causal et stable, doivent se trouver à l'intérieur du cercle unité [8].

1.4.6.4. Fonction de transfert d'un système régi par une équation aux différences :

Considérons le système régi par l'équation aux différences d'ordre N [8]:

$$\sum_{n=0}^N a_n \cdot y(k-n) = \sum_{m=0}^M b_m \cdot x(k-m) \quad (1.57)$$

où $y(k)$ est la réponse à l'excitation $x(k)$. En appliquant la transformée en Z aux deux membres de la relation, on aboutit à la fonction de transfert [8]:

$$H(z) = \frac{Y(z)}{X(z)} = \frac{\sum_{m=0}^M b_m \cdot z^{-m}}{\sum_{n=0}^N a_n \cdot z^{-n}} \quad (1.58)$$

1.5. SIGNAUX DISCRETS ALEATOIRES

1.5.1. Introduction:

Dans les paragraphes précédents, on a considéré les signaux déterministes. De tels signaux peuvent être représentés par la transformée en Z ou par la transformée de Fourier.

Une large classe de signaux physiques, par exemple les signaux

de communication, ne sont pas déterministes. Dans plusieurs situations, il est difficile, si ce n'est impossible, de donner une description précise de tels signaux. Ces signaux sont dits aléatoires. Leur représentation mathématique consiste en leur description en termes de moyennes. Plusieurs des propriétés de ces signaux sont déterministes et leur transformée en Z ou leur transformée de Fourier possède une interprétation particulière.

1.5.2. Processus aléatoires:

Un processus aléatoire est une famille indexée de variables aléatoires $\{X_n\}$. La famille de variables aléatoires est caractérisée par un ensemble de fonctions de distribution de probabilité qui est en général une fonction de l'indexe n dénotant le temps.

Une variable aléatoire individuelle X_n est décrite par la fonction de répartition de probabilité:

$$F(x_n, n) = \text{prob}(X_n \leq x_n) \quad (1.59)$$

où X_n dénote la variable aléatoire et x_n est une valeur particulière de X_n . Si X_n prend ses valeurs sur un ensemble continu, alors on peut obtenir la fonction densité de probabilité en dérivant la fonction de répartition:

$$f(x_n, n) = \frac{\partial F(x_n, n)}{\partial x_n} \quad (1.60)$$

$$\text{où} \quad F(x_n, n) = \int_{-\infty}^{x_n} f(x, n) dx \quad (1.61)$$

Si la variable aléatoire est quantifiée, alors elle prend des valeurs sur un ensemble dénombrable. Dans ce cas, la dérivée n'existe pas et on définit à sa place la fonction de masse de probabilité par :

$$f(x_n, n) = \text{Prob}(X_n = x_n) \quad (1.62)$$

Dans ce cas, la fonction de répartition de probabilité est donnée par :

$$F(x_n, n) = \text{Prob}(X_n \leq x_n) = \sum_{x \leq x_n} f(x, n) \quad (1.63)$$

La dépendance entre deux variables aléatoires X_n et X_m d'un

processus aléatoire est décrite par la fonction de répartition jointe (ou du deuxième ordre) :

$$F(x_n, x_m; n, m) = \text{Prob} [(X_n \leq x_n) \text{ et } (X_m \leq x_m)] \quad (1.64)$$

ou dans le cas de variables aléatoires continues, par la densité de probabilité jointe :

$$f(x_n, x_m; n, m) = \frac{\partial^2 F[x_n, x_m; n, m]}{\partial x_n \partial x_m} \quad (1.65)$$

Dans le cas de variables aléatoires quantifiées, la fonction de masse de probabilité est définie par :

$$f(x_n, x_m; n, m) = \text{Prob} [(X_n = x_n) \text{ et } (X_m = x_m)] \quad (1.66)$$

Une caractérisation complète d'un processus aléatoire exige la spécification de toutes les fonctions de répartition jointes. Dans le cas où toutes les fonctions de probabilité sont indépendantes de l'origine du temps, le processus aléatoire est dit stationnaire. Par exemple, la fonction de répartition du second ordre d'un processus stationnaire satisfait :

$$F(x_{n+k}, x_{m+k}; n+k, m+k) = F(x_n, x_m; n, m) \quad (1.67)$$

1.5.3. Moyennes:

Il est souvent utile de caractériser une variable aléatoire par des moyennes telles que l'espérance et la variance. Puisqu'un processus aléatoire est une famille indexée de variables aléatoires, on peut caractériser le processus par des moyennes statistiques des variables aléatoires comprenant le processus aléatoire. De telles moyennes sont appelées moyennes d'ensemble.

1.5.3.1 Définitions:

. L'espérance ou la moyenne d'un processus est définie par

$$m_x = E[X_n] = \int_{-\infty}^{+\infty} x \cdot f(x, n) dx \quad (1.68)$$

où E dénote l'espérance mathématique.

. La valeur quadratique moyenne de X_n est la moyenne de X_n^2

$$E [X_n^2] = \text{moyenne quadratique} = \int_{-\infty}^{+\infty} x^2 \cdot f(x, n) dx \quad (1.69)$$

Elle est parfois appelée la puissance moyenne.

. La variance de X_n est la valeur quadratique moyenne de $(X_n - m_{X_n})$

$$\begin{aligned} \text{Variance} &= E [(X_n - m_{X_n})^2] = \sigma_{X_n}^2 \\ &= E [X_n^2] - m_{X_n}^2 \\ &= \text{moyenne quadratique} - (\text{espérance})^2 \end{aligned} \quad (1.70)$$

. L'autocorrélation est définie par :

$$\begin{aligned} \phi_{xx}(n, m) &= E [X_n \cdot X_m^*] \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} x_n \cdot x_m^* f(x_n, x_m; n, m) dx_n \cdot dx_m \end{aligned} \quad (1.71)$$

. L'autocovariance d'un processus est définie par :

$$\begin{aligned} \gamma_{xx}(n, m) &= E [(X_n - m_{X_n}) (X_m - m_{X_m})^*] \\ &= \phi_{xx}(n, m) - m_{X_n} m_{X_m} \end{aligned} \quad (1.72)$$

. L'intercorrélation de deux processus $\{X_n\}$ et $\{Y_n\}$ est définie par:

$$\phi_{xy} = E [X_n Y_m^*] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} x \cdot y^* \cdot f(x, y; n, m) dx \cdot dy \quad (1.73)$$

. La fonction d'intercovariance est définie par :

$$\begin{aligned} \gamma_{xy} &= E [(X_n - m_{X_n}) (Y_m - m_{Y_m})^*] \\ &= \phi_{xy}(n, m) - m_{X_n} \cdot m_{Y_m} \end{aligned} \quad (1.74)$$

Cas des signaux stationnaires

Un processus stationnaire est caractérisé par le fait que ses propriétés statistiques ne dépendent pas de l'origine du temps. Ceci implique que la fonction de répartition du premier ordre est indépendante du temps et la fonction de répartition du second ordre satisfait l'équation (1.67). Il en suit que la répartition du second ordre dépend uniquement de la différence de temps $(m-n)$. L'espérance, la variance et la moyenne quadratique sont indépendantes du temps, et de ce fait, elles sont constantes. L'autocorrélation et l'autocovariance dépendent uniquement de la différence de temps $(m-n)$. Ceci est indiqué par les équations

suivantes :

$$m_X = E [X_n] = \text{constante} \quad (1.75)$$

$$\begin{aligned} \sigma_X^2 &= E [(X_n - m_X)^2] = E [X_n^2] - (m_X)^2 \\ &= \text{constante} \end{aligned} \quad (1.76)$$

$$\phi_{xx}(n, n+m) = \phi_{xx}(m) = E [X_n \cdot X_{n+m}^*] \quad (1.77)$$

$$\gamma_{xx}(n, n+m) = \phi_{xx}(m) - (m_X)^2 \quad (1.78)$$

1.5.3.2 Moyennes temporelles

On a vu que la notion d'un ensemble de signaux nous permet d'utiliser la théorie des probabilités dans la représentation des processus aléatoires. Cependant, en pratique, on préfère s'occuper d'une seule séquence plutôt que d'un ensemble de séquences. Par exemple, on désire déduire certaines moyennes des processus aléatoires à partir de mesures faites sur un seul membre de l'ensemble. Pour formaliser ces notions, on définit la moyenne temporelle du processus aléatoire par :

$$\langle X_n \rangle = \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{n=-N}^N X_n \quad (1.79)$$

L'autocorrélation temporelle est définie par :

$$\langle X_n X_{n+m} \rangle = \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{n=-N}^N X_n \cdot X_{n+m}^* \quad (1.80)$$

Dans le cas des processus ergodiques, les moyennes temporelles sont égales aux moyennes statistiques [27], ce qui donne :

$$\langle X_n \rangle = m_x \quad (1.81)$$

et $\langle X_n, X_{n+m} \rangle = \phi_{xx}(m)$

En pratique, on suppose qu'une séquence donnée est une séquence échantillon d'un processus aléatoire ergodique. Donc, les moyennes peuvent être calculées à partir d'une séquence unique. En général, on ne peut calculer les limites des équations (1.79) et (1.80), mais les quantités :

$$\langle X(n) \rangle_N = \frac{1}{2N+1} \sum_{n=-N}^N X(n) \quad (1.82)$$

et

$$\langle X(n) \cdot X(n+m) \rangle_N = \frac{1}{2N+1} \sum_{n=-N}^N X(n) \cdot X^*(n+m) \quad (1.83)$$

Ces quantités sont les estimations de la moyenne et de l'autocorrélation.

1.5.4. Représentation spectrale

1.5.4.1 Propriétés de la corrélation et de la covariance de séquence :

Considérons deux processus aléatoires réels stationnaires $\{X_n\}$ et $\{Y_n\}$ avec l'autocorrélation, l'autocovariance et l'intercovariance, données par [7]:

$$\phi_{xx}(m) = E [X_n X_{n+m}] \quad (1.84)$$

$$\gamma_{xx}(m) = E [(X_n - m_x)(X_{n+m} - m_x)] \quad (1.85)$$

$$\phi_{xy}(m) = E [X_n Y_{n+m}] \quad (1.86)$$

$$\gamma_{xy}(m) = E [(X_n - m_x)(Y_{n+m} - m_y)] \quad (1.87)$$

On a les propriétés suivantes :

1^{ère} Propriété :

$$\gamma_{xx}(m) = \phi_{xx}(m) - m_x^2 \quad (1.88.a)$$

$$\gamma_{xy}(m) = \phi_{xy}(m) - m_x m_y \quad (1.88.b)$$

2^{ème} Propriété :

$$\phi_{xx}(0) = E [X_n^2] : \text{valeur quadratique moyenne} \quad (1.89.a)$$

$$\gamma_{xx}(0) = \sigma_x^2 : \text{variance} \quad (1.89.b)$$

3^{ème} Propriété :

$$\phi_{xx}(m) = \phi_{xx}(-m) \quad (1.90.a)$$

$$\gamma_{xx}(m) = \gamma_{xx}(-m) \quad (1.90.b)$$

$$\phi_{xy}(m) = \phi_{yx}(-m) \quad (1.90.c)$$

$$\gamma_{xy}(m) = \gamma_{yx}(-m) \quad (1.90.d)$$

4^{ème} Propriété :

$$|\phi_{xx}(m)| \leq [\phi_{xx}(0) \cdot \phi_{yy}(0)]^{1/2} \quad (1.91.a)$$

$$|\gamma_{xy}(m)| \leq [\gamma_{xx}(0) \cdot \gamma_{yy}(0)]^{1/2} \quad (1.91.b)$$

En particulier, on a :

$$|\phi_{xx}(m)| \leq \phi_{xx}(0) \quad (1.92.a)$$

$$|\gamma_{xx}(m)| \leq \gamma_{xx}(0) \quad (1.92.b)$$

5^{ème} Propriété : Si $Y_n = X_{n-n_0}$, alors

$$\phi_{yy}(m) = \phi_{xx}(m) \quad (1.93.a)$$

$$\gamma_{yy}(m) = \gamma_{xx}(m) \quad (1.93.b)$$

6^{ème} Propriété :

$$\lim_{m \rightarrow \infty} \phi_{xx}(m) = (E[X_n])^2 = m_x^2 \quad (1.94.a)$$

$$\lim_{m \rightarrow \infty} \gamma_{xx}(m) = 0 \quad (1.94.b)$$

$$\lim_{m \rightarrow \infty} \phi_{xy}(m) = m_x \cdot m_y \quad (1.94.c)$$

$$\lim_{m \rightarrow \infty} \gamma_{xy}(m) = 0 \quad (1.94.d)$$

1.5.4.2. Transformée en Z

Soient $\phi_{xx}(z)$, $\Gamma_{xx}(z)$, $\phi_{xy}(z)$ et $\Gamma_{xy}(z)$ les transformées en z de $\phi_{xx}(m)$, $\gamma_{xx}(m)$, $\phi_{xy}(m)$ et $\gamma_{xy}(m)$ respectivement.

Des équations (1.94.a) et (1.94.c), on note immédiatement que les transformées en Z de $\phi_{xx}(m)$ et $\phi_{xy}(m)$ n'existent que lorsque $m_x = 0$. Dans ce cas, $\phi_{xx}(z) = \Gamma_{xx}(z)$ et $\phi_{xy}(z) = \Gamma_{xy}(z)$.

Les propriétés de la transformée en Z sont données comme suit :

1^{ère} Propriété :

$$\sigma_x^2 = \frac{1}{2\pi j} \oint_C \Gamma_{xx}(z) z^{-1} dz \quad (1.95)$$

où c est un contour fermé dans la région de convergence de $\Gamma_{xx}(z)$

2^{ème} Propriété :

$$\Gamma_{xx}(z) = \Gamma_{xx}(1/z) \quad (1.96.a)$$

$$\Gamma_{xy}(z) = \Gamma_{yx}^*(1/z^*) \quad (1.96.b)$$

Les équations (1.96) découlent directement de la propriété 3 de la section 1.5.4.1. Comme conséquence, la région de convergence de

$\Gamma_{xx}(z)$ doit être de la forme :

$$R_a < |z| < 1/R_a$$

En plus, en vertu de l'équation (1.94.b), la région de convergence doit contenir le cercle unité, i.e $0 < R_a < 1$

1.5.4.3 Spectre de puissance :

Puisque la région de convergence contient le cercle unité, on peut exprimer l'équation (1.95) par :

$$\sigma_x^2 = \frac{1}{2\pi} \int_{-\pi}^{+\pi} P_{xx}(\omega) d\omega \quad (1.97)$$

où on définit

$$P_{xx}(\omega) = \Gamma_{xx}(e^{j\omega}) \quad (1.98)$$

Quand $m_x = 0$, la variance est égale à la puissance moyenne. Ainsi l'aire au dessous de $P_{xx}(\omega)$ pour $-\pi < \omega < \pi$ est proportionnelle à la puissance moyenne du signal. De ce fait, l'intégrale de $P_{xx}(\omega)$ sur une bande de fréquence est proportionnelle à la puissance du signal dans cette bande. Pour ces raisons, $P_{xx}(\omega)$ est appelée Densité Spectrale de Puissance. Et il est commun de définir la densité Spectrale comme la transformée de Fourier de l'autocorrélation plutôt que de l'autocovariance.

De la propriété 2 de la section (1.5.4.2), on obtient :

$$P_{xx}(\omega) = P_{xx}(-\omega) \quad (1.99)$$

On définit similairement, la densité spectrale d'interpuissance par :

$$P_{xy}(\omega) = \Gamma_{xy}(e^{j\omega}) \quad (1.100)$$

De la propriété 2 de la section (1.5.4.2), on obtient :

$$P_{xy}(\omega) = P_{xy}^*(-\omega)$$

1.5.5 Réponse des systèmes linéaires aux signaux discrets :

Considérons un système linéaire invariant et stable ayant la réponse impulsionnelle $h(n)$. Et soit $x(n)$ une séquence d'un

processus aléatoire stationnaire au sens large. On excite le système avec la séquence $x(n)$. La réponse du système est une fonction échantillon d'un processus aléatoire lié au processus d'entrée par la somme de convolution [7]:

$$y(n) = \sum_{k=-\infty}^{+\infty} h(k) \cdot x(n-k) \quad (1.101)$$

On désire tirer les caractéristiques du processus de sortie à partir de celle du processus d'entrée. Celui-ci ayant une moyenne m_x . Une fonction d'autocorrélation $\phi_{xx}(m)$ et une variance σ_x^2 .

- La moyenne du processus de sortie est :

$$m_y = E [Y(n)] = \sum_{k=-\infty}^{+\infty} h(k) \cdot E [x(n-k)] \quad (1.102)$$

$$= m_x \sum_{k=-\infty}^{+\infty} h(k) \quad (1.103)$$

ou en termes de la réponse fréquentielle du système :

$$m_y = H(0) \cdot m_x \quad (1.104)$$

Puisque l'entrée est stationnaire, on voit que la moyenne de la sortie est constante.

- L'autocorrélation du processus de sortie est [7]:

$$\phi_{yy}(m) = \sum_{l=-\infty}^{+\infty} \phi_{xx}(m-l) \sum_{k=-\infty}^{+\infty} h(k) \cdot h(l+k) \quad (1.105)$$

$$= \sum_{l=-\infty}^{+\infty} \phi_{xx}(m-l) \cdot v(l)$$

où on définit :

$$v(l) = \sum_{k=-\infty}^{+\infty} h(k) \cdot h(l+k) \quad (1.106)$$

$v(l)$ est appelée fonction d'autocorrélation de $h(n)$, qui est aussi la convolution de $h(n)$ avec $h(-n)$. Ceci nous permet d'écrire l'équation (1.105) sous la forme [27]:

$$\phi_{yy}(m) = \phi_{xx}(m) * h(m) * h(-m) \quad (1.107)$$

- L'intercorrélation entre l'entrée et la sortie est donnée par [6]:

$$\phi_{xy}(m) = \sum_{k=-\infty}^{+\infty} h(k) \cdot \phi_{xx}(m-k) \quad (1.108)$$

ou par la convolution [27]:

$$\phi_{xy}(m) = h(m) * \phi_{xx}(m) \quad (1.109)$$

Si on suppose que $m_x = 0$, la transformée en Z de la fonction d'autocorrélacion existe. Alors de l'équation (1.107), on tire [6]:

$$\phi_{yy}(z) = H(z) \cdot H(z^{-1}) \cdot \phi_{xx}(z) \quad (1.110)$$

- La densité spectrale de puissance est :

$$P_{yy}(\omega) = |H(\omega)|^2 \cdot P_{xx}(\omega) \quad (1.111)$$

- La puissance moyenne totale est [7]

$$\begin{aligned} \sigma_y^2 = \phi_{yy}(0) &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} P_{xx}(\omega) \cdot d\omega \\ &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} |H(\omega)|^2 \cdot P_{xx}(\omega) \cdot d\omega \end{aligned} \quad (1.112)$$

De l'équation (1.109), on obtient :

$$\phi_{xy}(z) = H(z) \cdot \phi_{xx}(z) \quad (1.113)$$

. La densité spectrale d'interpuissance est [6]

$$P_{xy}(\omega) = H(\omega) \cdot P_{xx}(\omega) \quad (1.114)$$

1.5.6. Bruit blanc:

Le bruit blanc est un processus aléatoire dont le spectre de puissance est constant dans toute l'étendue de la fréquence [6].

De cette définition, la fonction d'autocorrélacion et la densité spectrale de puissance sont données par [6]:

$$\phi_{xx}(m) = \sigma_x^2 \delta(m) \quad (1.115)$$

et $P_{xx}(\omega) = \phi_{xx}(e^{j\omega}) = \sigma_x^2$

Si on excite le système de la section (1.5.5) par un bruit blanc, les résultats énoncés deviennent comme suit :

- La fonction d'autocorrélation (équation 1.107) devient :

$$\phi_{yy}(m) = \sigma_x^2 (h(m) * h(-m)) \quad (1.116)$$

- Sa transformée en Z (équation 1.110) devient :

$$\phi_{yy}(z) = \sigma_x^2 H(z) \cdot H(z^{-1}) \quad (1.117)$$

- La densité spectrale (équation 1.111) devient:

$$P_{yy}(\omega) = \sigma_x^2 |H(\omega)|^2 \quad (1.118)$$

- La puissance moyenne totale (équation 1.112) est donnée par :

$$\sigma_y^2 = \frac{\sigma_x^2}{2\pi} \int_{-\pi}^{+\pi} |H(\omega)|^2 d\omega \quad (1.119)$$

ou par $\sigma_y^2 = \sigma_x^2 \sum_{n=0}^{\infty} h^2(n)$ (1.120)

- L'intercorrélation (équation 1.113) devient :

$$\phi_{xy}(m) = \sigma_x^2 \cdot h(m) \quad (1.121)$$

- Sa transformée en Z devient :

$$\phi_{xy}(z) = \sigma_x^2 \cdot \phi_{xx}(z) \quad (1.122)$$

- La densité spectrale d'interpuissance (équation 1.114) devient [7]:

$$P_{xy}(\omega) = \sigma_x^2 \cdot H(\omega) \quad (1.123)$$

DISCRETE

La transformation de Fourier discrète (TFD) est l'une des opérations les plus importantes du traitement numérique du signal. En plus de son aspect théorique, la TFD joue un rôle central dans l'implantation d'une variété d'algorithmes de traitement numérique du signal. Ce rôle est dû à l'existence d'un algorithme efficace pour le calcul de la TFD [12]. Cet algorithme est la transformation de Fourier rapide (FFT) qu'on présentera au paragraphe (2.7)

Avant d'aborder la TFD, on traitera la représentation des séquences périodiques, ou la série de Fourier discrète (SFD).

2.1 Série de Fourier Discrète:

Soit un signal discret $x_p(n)$ périodique de période N . Il est possible de représenter $x_p(n)$ en terme d'une série de Fourier par la somme de séquences de sinus et de cosinus, ou par des séquences exponentielles complexes, avec des fréquences qui sont des multiples entiers de la fréquence fondamentale $1/N$ associée à $x_p(n)$. Mais en opposition avec les séries de Fourier pour les fonctions continues périodiques, il n'y a que N exponentielles complexes ayant une fréquence qui est un multiple entier de la fréquence fondamentale $1/N$. Ainsi la série de Fourier d'une séquence périodique $x(n)$ ne contient que N exponentielles complexes [7]:

$$x_p(n) = \frac{1}{N} \sum_{k=0}^{N-1} X_p(k) e^{j(2\pi/N)nk} \quad (2.1)$$

Les coefficients $X(k)$ sont donnés par la relation [7] :

$$X_p(k) = \sum_{n=0}^{N-1} x_p(n) e^{-j(2\pi/N)nk} \quad (2.2)$$

De cette relation, on remarque que la séquence donnant $X_p(k)$ est périodique de période N . Par convenance, on écrit les relations (2.1) et (2.2) en fonction de W_N défini par:

$$W_N = e^{-j(2\pi/N)} \quad (2.3)$$

Ceci donne :

$$X_p(k) = \sum_{n=0}^{N-1} x_p(n) W_n^{nk} \quad (2.4)$$

$$x_p(n) = \frac{1}{N} \sum_{k=0}^{N-1} X_p(k) W_n^{-nk} \quad (2.5)$$

Les relations (2.4) et (2.5) sont appelées respectivement analyse et synthèse de la série de Fourier. Ainsi une séquence périodique est complètement représentée par sa série de Fourier.

2.2. Propriétés de la série de Fourier :

2.2.1 Linéarité :

Si deux séquences périodiques $x_{p_1}(n)$ et $x_{p_2}(n)$, de périodes égales à N sont combinées pour former la séquence $x_{p_3}(n)$ telle que:

$$x_{p_3}(n) = a x_{p_1}(n) + b x_{p_2}(n)$$

alors les coefficients de Fourier de $x_3(n)$ sont donnés par [7]

$$X_{p_3}(k) = a X_{p_1}(k) + b X_{p_2}(k) \quad (2.6)$$

où toutes les séquences sont de période N , et a et b sont des scalaires.

2.2.2 Décalage d'une séquence :

Si une séquence périodique $x_p(n)$ a les coefficients de Fourier $X_p(k)$, alors la séquence $x_p(n+m)$ a les coefficients de Fourier $W_n^{-mk} \cdot X_p(k)$. m étant un entier.

2.2.3 Convolution périodique :

Soient $x_{p_1}(n)$ et $x_{p_2}(n)$ deux séquences périodiques de période N ayant les coefficients de Fourier respectifs $X_{p_1}(k)$ et $X_{p_2}(k)$. On forme la séquence $x_{p_3}(n)$ tel que l'on ait [7]:

$$x_{p_3}(n) = \sum_{m=0}^{N-1} x_{p_1}(m) x_{p_2}(n-m) \quad (2.7)$$

Les coefficients de Fourier de la séquence $x_3(n)$ sont donnés par [7]:

$$X_{p_3}(k) = X_{p_1}(k) \cdot X_{p_2}(k) \quad (2.8)$$

La relation (2.7) est appelée convolution circulaire. Elle présente deux aspects qui la diffèrent de la convolution de

séquences apériodiques :

* D'une part, $x_{p1}(n)$ et $x_{p2}(n)$ sont périodiques, de période N et ainsi est leur produit de convolution $x_{p3}(n)$.

* D'autre part, la sommation est évaluée sur une période.

L'illustration de ce type de convolution est faite par la figure 2.1. Ce qui est important, c'est que quand une période glisse vers l'extérieur de l'intervalle de sommation, l'autre période y pénètre. L'oubli de ceci peut conduire à de graves erreurs dans l'interprétation des résultats.

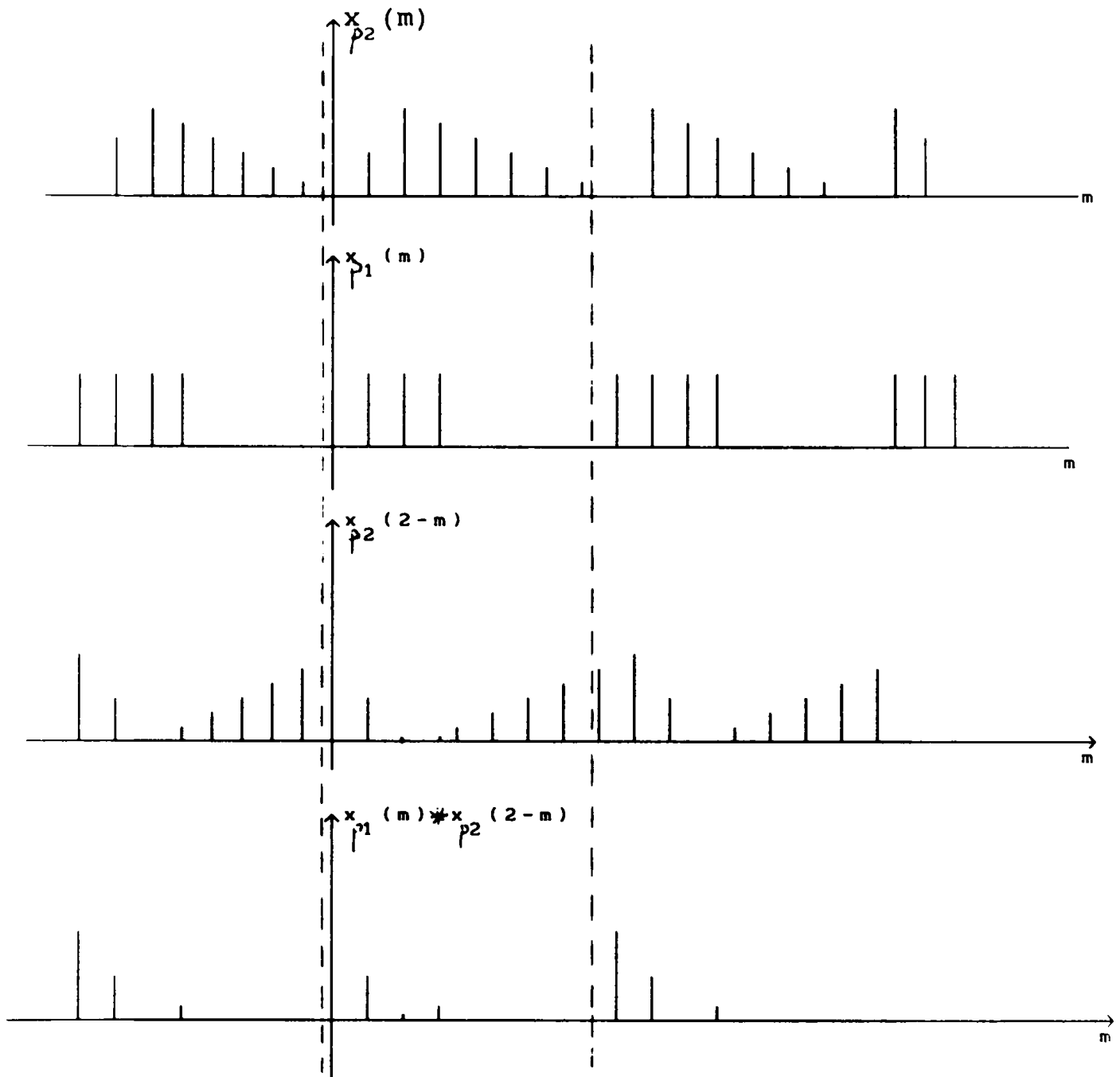


Figure 2.1 : Procédure de formation de la convolution périodique de deux séquences périodiques.

2.3 Transformation de Fourier Discrète :

Avec une interprétation correcte, la représentation du paragraphe précédent peut être appliquée aux signaux à durée finie (limitée). La représentation de Fourier résultante est appelée transformation de Fourier Discrète (TFD).

Une séquence à durée finie de longueur N peut être représentée par une séquence périodique, de période N , dont chaque période est identique à la séquence de durée finie. Du fait que la séquence périodique possède une décomposition unique en série de Fourier, il en sera de même pour la séquence originale à durée finie, puisqu'on peut calculer une période unique d'une séquence périodique à partir de sa série de Fourier Discrète.

Soit $x(n)$ un signal à durée finie de longueur N , de sorte que $x(n)$ soit nul à l'extérieur de l'intervalle $[0, N-1]$. La séquence périodique correspondante $x_p(n)$ est donnée par [7] :

$$x_p(n) = \sum_{r=-\infty}^{+\infty} x(n + rN) \quad (2.9)$$

La séquence $x(n)$ est obtenue en extrayant une période de $x_p(n)$ [7] :

$$x(n) = \begin{cases} x_p(n) & , \quad 0 \leq n \leq N-1 \\ 0 & , \quad \text{ailleurs} \end{cases} \quad (2.10)$$

Une relation similaire est donnée pour les coefficients $X(k)$:

$$X(k) = \begin{cases} X_p(k) & , \quad 0 \leq k \leq N-1 \\ 0 & , \quad \text{ailleurs} \end{cases} \quad (2.11)$$

Puisque les relations (2.4) et (2.5) nécessitent uniquement l'intervalle entre 0 et $N-1$, les relations (2.10) et (2.11) deviennent :

$$X(k) = \sum_{n=0}^{N-1} x(n) W_N^{nk} \quad 0 \leq k \leq N-1 \quad (2.12)$$

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k) W_N^{-nk} \quad 0 \leq n \leq N-1 \quad (2.13)$$

Les relations (2.12) et (2.13) sont appelées transformations de Fourier Discrète (TFD) avec la relation (2.12) représentant l'analyse de la TFD et la relation (2.13) représentant la synthèse

de la TFD ou la TFD inverse.

Là où on applique la TFD, on doit se rappeler qu'une séquence à durée finie est représentée comme une période d'une séquence périodique à laquelle on applique la SFD. L'oubli de ceci peut conduire à des erreurs.

2.4 Propriétés de la TFD :

2.4.1 Linéarité :

Si deux séquences de durée finie $x_1(n)$ et $x_2(n)$ sont combinées pour former la séquences $x_3(n)$, telle que :

$$x_3(n) = ax_1(n) + bx_2(n)$$

alors la TFD de $x_3(n)$ est :

$$X_3(k) = a X_1(k) + b X_2(k) \quad (2.14)$$

Si $x_1(n)$ est de durée N_1 et $x_2(n)$ est de durée N_2 , alors $x_3(n)$ est de durée N_3 tel que :

$$N_3 = \max [N_1, N_2]$$

Ainsi les TFD doivent être calculées avec $N = N_3$. Si par exemple N_1 est inférieur à N_2 , alors $X_1(k)$ est la TFD de la séquence $x_1(n)$ prolongée par $N_2 - N_1$ zéros.

2.4.2 Décalage circulaire d'une séquence :

Considérons un signal à durée limitée $x(n)$, sa version périodique $x_p(n)$ et sa version périodique décalée $x_p(n+m)$ représentés sur les figures (2.2.a) à (2.2.c) respectivement. La séquence de durée finie, notée $x_1(n)$, obtenue en extrayant une période de la séquence $x_p(n+m)$ est montrée sur la figure (2.2.d). La comparaison entre les figures (2.9.a) et (2.9.b) indique clairement que $x_1(n)$ ne correspond pas à un décalage linéaire de $x(n)$.

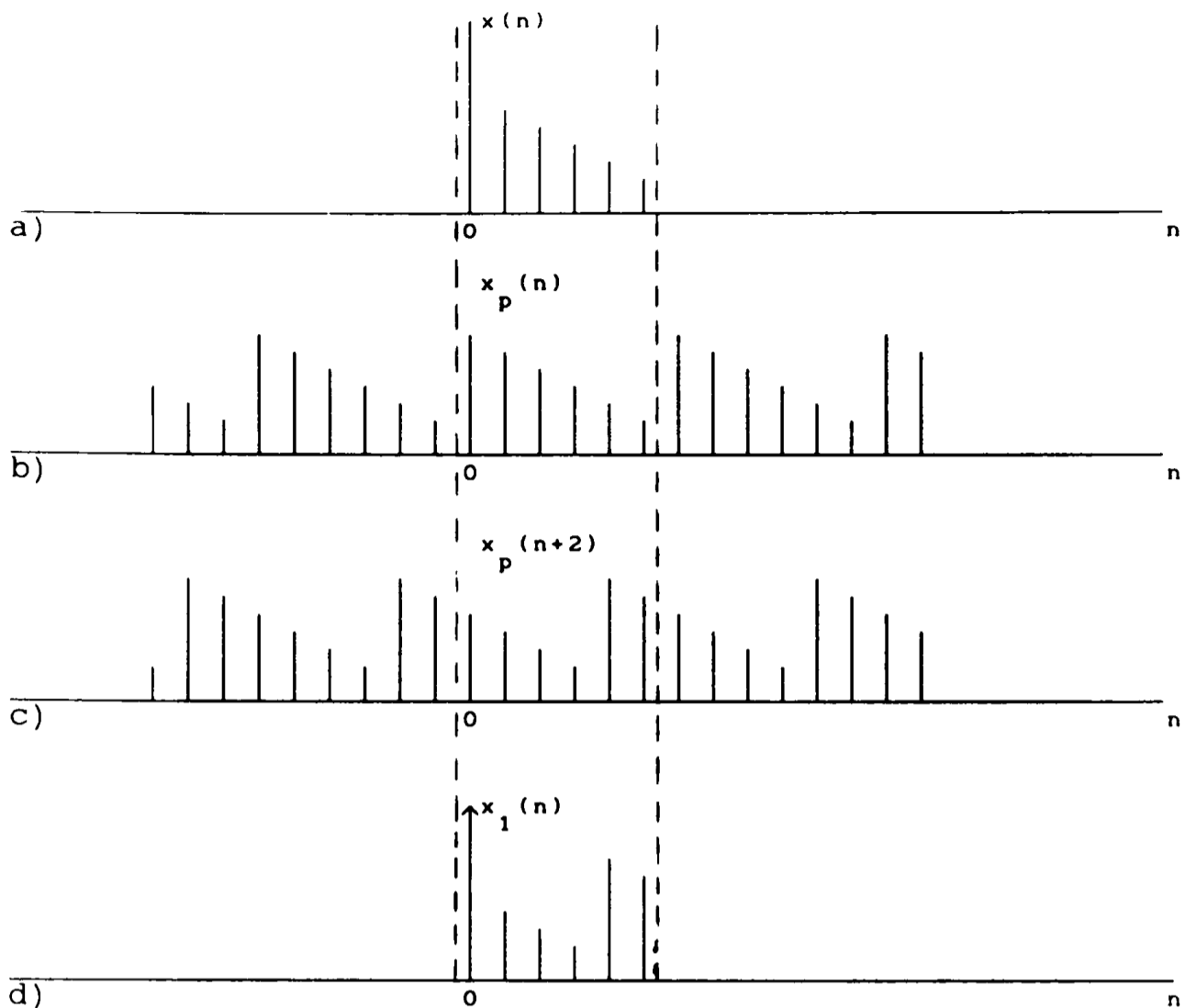


Figure 2.2 : Decalage circulaire d'une sequence

De ce fait, la TFD de $x_1(n)$, donnée par: $X_1(k) = W_N^{-mk} X(k)$ ne correspond pas à la TFD de $x(n+m)$.

2.4.3 Propriété de symétrie :

Les principales propriétés de symétrie sont :

- Si la TFD de la séquence $x(n)$ est $X(k)$, alors la TFD de la séquence $x^*(n)$ est $X^*(-k)$ avec $*$ indiquant l'expression conjuguée.

- Si $x(n)$ est réelle, on a:

$\text{Re}[X(k)]$ est une fonction paire

et $\text{Im}[X(k)]$ est une fonction impaire

où $\text{Re}[\cdot]$ dénote la partie réelle et $\text{Im}[\cdot]$ dénote la partie imaginaire.

2.4.4 Convolution circulaire :

Soient $x_1(n)$ et $x_2(n)$ deux séquences à durée finie N , ayant les TFD respectives $X_1(k)$ et $X_2(k)$. La séquence $x_3(n)$ de durée N ,

correspondant à $X_1(k) \cdot X_2(k)$ est obtenue par extraction d'une période de la séquence $x_p(n)$ correspondant à la convolution circulaire des séquences $x_{p1}(n)$ et $x_{p2}(n)$ les versions périodiques des séquences $x_1(n)$ et $x_2(n)$. Ceci diffère de la convolution linéaire par le fait que la séquence obtenue par la convolution de $x_1(n)$ et $x_2(n)$ donne naissance à une séquence de durée $2N-1$, et que lors de la sommation dans la convolution linéaire, les décalages sont linéaires alors que dans la convolution circulaire, les décalages sont circulaires. Ce qui provoque un chevauchement indésirable (Figure 2.1).

Pour remédier à ce problème, on prolonge les séquences $x_1(n)$ et $x_2(n)$ par $(N-1)$ zéros et on calcule ainsi leur DFT sur $(2N-1)$ points (figure 2.3).

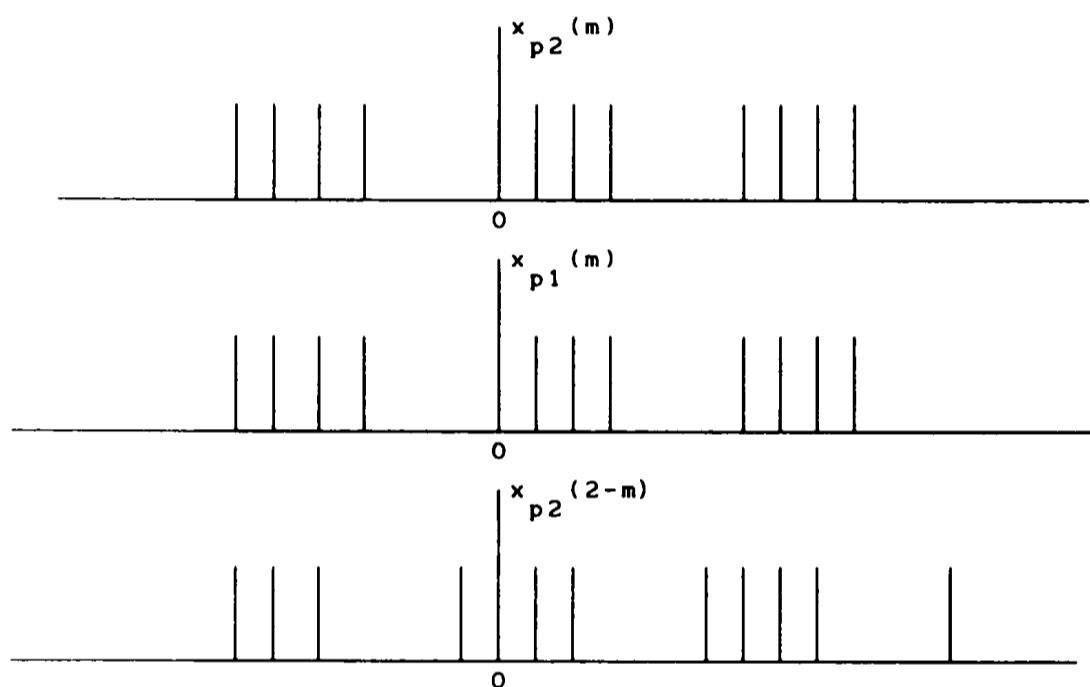


Figure 2.3 : Séquences périodiques obtenues en prolongeant les séquences originales par des zéros.

3.5 Transformation de Fourier Discrète pour les signaux à durée illimitée.

On a défini la TFD pour les signaux à durée limitée. Pour les signaux à durée illimitée, elle ne peut être définie qu'approximativement en limitant la durée du signal par un moyen approprié.

La limitation de la durée du signal est obtenue en le multipliant par une fenêtre de troncature. Parmi les fenêtres qu'on utilise, on peut citer les fenêtres de Hamming, Hanning, Bartlett, Kaiser et la fenêtre rectangulaire. La définition et les caractéristiques de ces fenêtres peuvent être trouvées dans la référence [13].

2.6 Convolution linéaire :

L'existence d'un algorithme efficace pour le calcul de la TFD de séquence à durée limitée rend efficace l'implantation de la convolution linéaire de deux séquences en calculant la TFD inverse du produit de leur TFD.

Si on a deux séquences $x_1(n)$ et $x_2(n)$ à durée limitée N , alors la séquence $x_3(n)$ obtenue par la convolution linéaire de $x_1(n)$ et de $x_2(n)$ est à durée limitée $(2N-1)$ on a :

$$x_3(n) = \sum_{m=0}^{N-1} x_1(m) \cdot x_2(n-m) \quad (2.15)$$

Les TFD de $x_1(n)$ et de $x_2(n)$ doivent être calculées de façon à assurer l'obtention d'une convolution linéaire. Pour cette raison les TFD sont calculées sur $(2N-1)$ points avec les séquences $x_1(n)$ et $x_2(n)$ prolongées de $(N-1)$ zéros. On a :

$$X_1(k) = \sum_{n=0}^{2N-1} x_1(n) \cdot W_{2N-1}^{nk}$$
$$X_2(k) = \sum_{n=0}^{2N-1} x_2(n) \cdot W_{2N-1}^{nk} \quad (2.16)$$

$$x_3(n) = \frac{1}{2N-1} \sum_{k=0}^{2N-1} [X_1(k) \cdot X_2(k)] \cdot W_{2N-1}^{-nk}$$

Comme il sera montré au paragraphe suivant, une TFD sur N points nécessite $(N/2 \cdot \log_2(N))$ multiplications complexes. Sur cette base, on déduit que le calcul de $X_1(k)$, $X_2(k)$ et $x_3(n)$ exige $((2N-1)/2 \cdot \log_2(2N-1))$ multiplications complexes. En plus des $2N-1$ multiplications nécessaires pour obtenir le produit $X_1(k) \cdot X_2(k)$ ceci donne un total de $(3/2(2N-1)\log_2(2N-1) + 2N-1)$ multiplications complexes. Soit, lorsque N est assez grand un total de $(3N \cdot \log_2(N) + 5N)$ multiplications complexes. La relation (2.15) nécessite N^2 multiplications.

A titre de comparaison, on dresse la table 2.1 :

N	Eq 2.15	Eq 2.16
16	256	272
32	1024	640
64	4096	1472
1024	1.048576	35840

Table 2.1 : Nombre de multiplications nécessaires pour la convolution linéaire pour différentes valeurs de N

Cette table met en évidence l'avantage de l'implantation de la convolution linéaire par TFD surtout lorsque N est assez grand.

2.7 Transformation de Fourier Rapide FFT

Le calcul de la TFD selon la relation (2.12) exige N^2 multiplications complexes et (N^2-N) additions complexes. En appliquant la périodicité et la symétrie de l'exponentielle complexe W_N^{nk} , on peut diminuer considérablement le nombre de multiplications et d'additions nécessaires. On considère N une puissance de 2, $N = 2^M$. On a :

$$X(k) = \sum_{n=0}^{N-1} x(n) \cdot W_N^{nk}$$

En séparant les échantillons d'indices pairs et ceux d'indices impairs, on obtient [6]:

$$X(k) = \sum_{n=0}^{N/2-1} x(2n) \cdot W_N^{2nk} + \sum_{m=0}^{N/2-1} x(2m+1) \cdot W_N^{(2m+1)k} \quad (2.17)$$

avec $W_N^{2nk} = W_{N/2}^{nk}$

on obtient :

$$X(k) = \sum_{m=0}^{N/2-1} x(2m) \cdot W_{N/2}^{mk} + W_N^k \sum_{m=0}^{N/2-1} x(2m+1) \cdot W_{N/2}^{mk} \quad (2.18)$$

Chaque somme est une TFD sur N/2 points, avec :

$$X_1(k) = \sum_{m=0}^{N/2-1} x(2m) \cdot W_{N/2}^{mk} = \sum_{m=0}^{N/2-1} x_1(m) \cdot W_{N/2}^{mk}$$

et $X_2(k) = \sum_{m=0}^{N/2-1} x(2m+1) \cdot W_{N/2}^{mk} = \sum_{m=0}^{N/2-1} x_2(m) \cdot W_{N/2}^{mk}$

On obtient :

$$X(k) = X_1(k) + W_N^k \cdot X_2(k) , \quad k = 0, \dots, N/2-1 \quad (2.19)$$

$X_1(k)$ et $X_2(k)$ sont périodiques et de période $N/2$, les échantillons de $X(k)$ pour k compris entre $N/2$ et $N-1$ sont obtenus par la relation suivante :

$$\begin{aligned} X(k + \frac{N}{2}) &= X_1(k) - W_N^{\frac{N}{2}+k} X_2(k) \\ &= X_1(k) - W_N^k \cdot X_2(k) \end{aligned} \quad (2.20)$$

Puisque $W_N^{k+\frac{N}{2}} = -W_N^k$

Cette première décomposition est illustrée par la figure (2.4). Les équations (2.19) et (2.20) peuvent être rassemblées pour donner :

$$\begin{aligned} X(k) &= X_1(k) + W_N^k X_2(k) \\ X(k + \frac{N}{2}) &= X_1(k) - W_N^k X_2(k) \end{aligned} \quad (2.21)$$

Pour obtenir $X(k)$ à partir de la relation (2.21), il faut $N/2$ multiplications et N additions.

Pour obtenir $X_1(k)$, il faut $(N/2)^2$ multiplications et $[(\frac{N}{2})^2 - \frac{N}{2}]$ additions. La même chose est nécessaire pour obtenir $X_2(k)$.

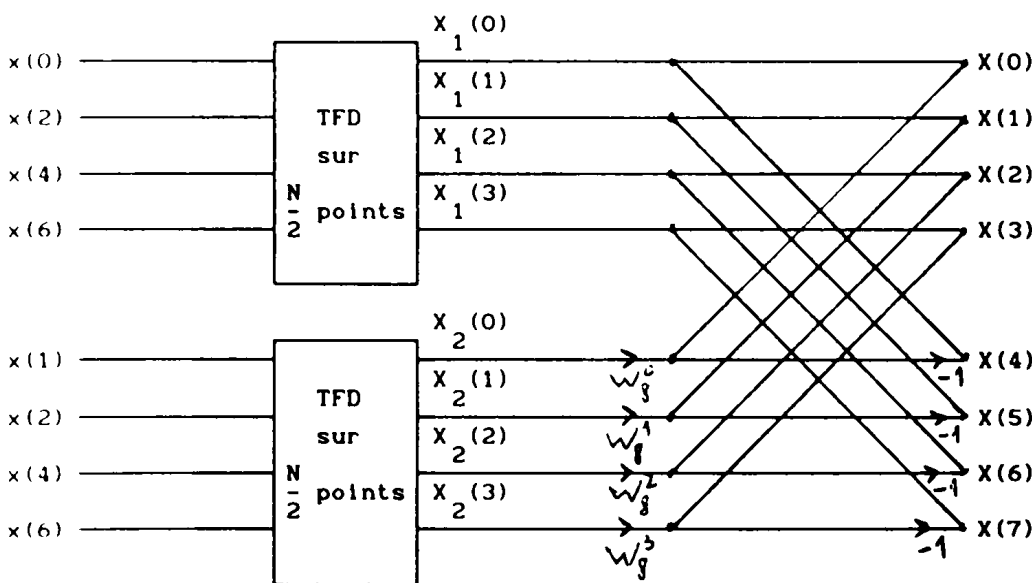


Figure 2.4 : Première décomposition dans le développement d'une TFD sur 8 points.

Soit un total de $(\frac{N}{2} + \frac{N}{2})$ multiplications et $(\frac{N}{2})$ additions. Le nombre d'opérations nécessaires a diminué presque de moitié.

On applique la même approche pour $X_1(k)$ et $X_2(k)$.

$$\begin{aligned} \text{On a : } X_1(k) &= \sum_{m=0}^{N/2-1} x_1(m) \cdot W_{N/2}^{nk} \\ &= \sum_{m=0}^{N/4-1} x_1(2m) \cdot W_{N/4}^{nk} + W_{N/2}^k \sum_{m=0}^{N/4-1} x_1(2m+1) \cdot W_{N/4}^{nk} \end{aligned}$$

Chaque somme est une TFD sur $N/4$ points. Ceci donne :

$$X_1(k) = X_{11}(k) + W_{N/2}^k \cdot X_{12}(k)$$

$X_{11}(k)$ et $X_{12}(k)$ sont périodiques et de période $N/2$, ainsi, on obtient :

$$X_1(k) = X_{11}(k) + W_{N/2}^k X_{12}(k)$$

$$X_1(k + \frac{N}{2}) = X_{11}(k) - W_{N/2}^k X_{12}(k)$$

La même chose peut être obtenue pour $X_2(k)$

$$X_2(k) = X_{21}(k) + W_{N/2}^k X_{12}(k)$$

$$X_2(k + \frac{N}{2}) = X_{21}(k) - W_{N/2}^k X_{12}(k)$$

Cette deuxième décomposition est illustrée par la figure (2.5).

La décomposition est répétée jusqu'à ce qu'une TFD de 2 points soit générée. Chaque décomposition est appelée étage. Le nombre d'étage est :

$$M = \log_2(N) \quad (2.22)$$

La décomposition complète de la TFD sur 8 points est illustrée par la figure (2.6).

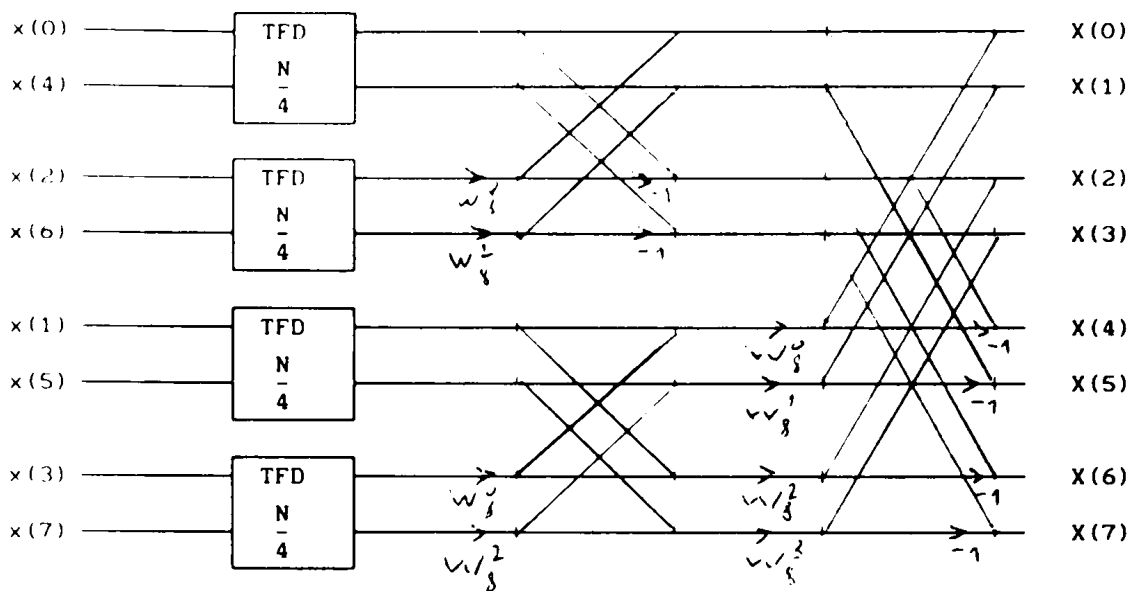


Figure 2.5 : Deuxieme decomposition de la TFD sur 8 points

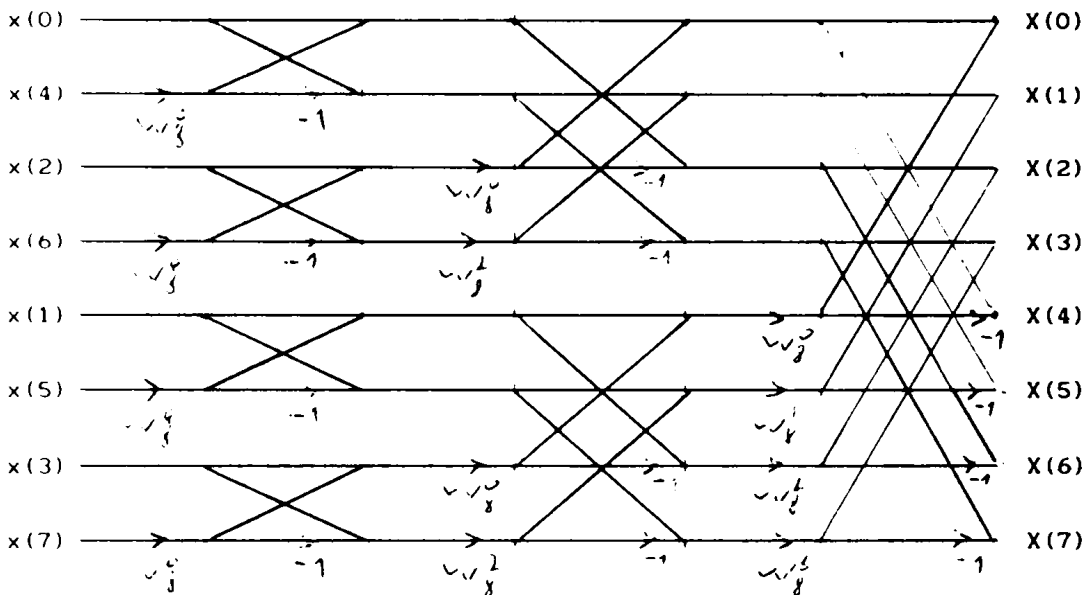


Figure 2.6 : Decomposition complete d'une TFD sur 8 points

L'équation (2.21) est communément appelée papillon. Son graphe de fluence illustré par la figure (2.7)

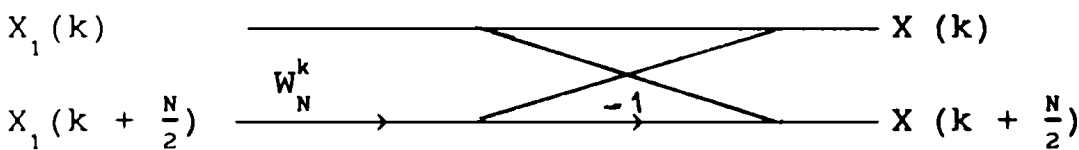


Figure 2.7 : Graphe de fluence d'un papillon

Le calcul d'un papillon nécessite une multiplication et deux additions complexes.

La décomposition complète de la TFD nécessite M étages et chaque étage nécessite N/2 papillons. Ce qui donne $(M \cdot \frac{N}{2})$ multiplications

et $(M.N)$ additions complexes. Soit :

$$\left(\frac{N}{2} \log_2 N\right) \text{ multiplications et } (N \log_2 N) \text{ additions}$$

comparé avec les (N^2) multiplications et (N^1-N) additions nécessaires pour le calcul direct, on voit l'avantage de la FFT.

La table 2.2 illustre une liste comparative :

N	Multiplication		Addition	
	Calcul direct	FFT	Calcul direct	FFT
16	256	32	240	64
64	4096	192	4032	384
512	262144	2304	261632	4608
1024	1048576	5120	1047552	10240

Table 2.2 : Table comparative des opérations nécessaires pour le calcul de la TFD

La décomposition présentée correspond à l'algorithme de la FFT le plus utilisé qu'on appelle FFT à entrelacement temporel (de l'anglais Decimation in Time FFT algorithm). Un autre algorithme de décomposition moins utilisé correspond à la FFT à entrelacement fréquentiel [15]. Il ne sera pas présenté ici.

2.7.1. Certains aspects de l'algorithme de la FFT :

- Une FFT est composée de M étages avec $M = \log_2 N$. Chaque étage contient $N/2$ papillons et contient des groupes.

- Tous les groupes d'un même étage ont le même ensemble d'exponentielles complexes W_N^k .

- Si on indexe les étages en fonction de leur position dans la décomposition dans le temps par m avec $m = 1 \dots M$, on tire les résultats suivants :

* Le nombre de groupes dans l'étage m est $N/2^m$

* Le nombre de papillons par groupe dans l'étage m est $(2^m - 1)$

* L'argument k des exponentielles complexes W_N^k d'un même groupe

dans l'étage m portent les valeurs suivantes :

$$k = \frac{iN}{2^m} \quad \text{avec } i = 0, \dots, 2^m - 1$$

- Les calculs sont fait sur place, c'est-à-dire que lors des

calculs intermédiaires, le même vecteur (tableau) qui contient la séquence d'entrée contiendra la séquence de sortie.

- A l'entrée de l'algorithme, les données doivent être réordonnées dans l'ordre bit inversé. Chaque indice d'un échantillon d'entrée est converti dans sa représentation en binaire, ensuite on lit cette représentation dans un ordre inverse ce qui donne l'indice de la nouvelle position de cet échantillon à l'entrée. Si on a $M = \log_2 N$ alors l'indice i de l'échantillon d'entrée s'écrit en binaire :

$$i = \sum_{r=0}^{M-1} B_r \cdot 2^r \quad \text{avec les } B_r = 0,1$$

L'index bit-inversé correspondant est :

$$i' = \sum_{r=0}^{M-1} B_r \cdot 2^{M-1-r} \quad B_r = 0,1$$

La figure 2.8 illustre l'opération d'inversion de bit pour $N=8$.

index de l'échantillon d'entrée		index bit-inverse de l'échantillon	
Decimal	Binaire	Binaire	Decimal
0	0 0 0	0 0 0	0
1	0 0 1	1 0 0	4
2	0 1 0	0 1 0	2
3	0 1 1	1 1 0	6
4	1 0 0	0 0 1	1
5	1 0 1	1 0 1	5
6	1 1 0	0 1 1	3
7	1 1 1	1 1 1	7

Figure 2.8 : Processus d'inversion de bit pour $N = 8$

- Le nombre d'exponentielles complexes nécessaires pour la FFT est $N/2$

Ainsi l'algorithme de la FFT se divise en deux étapes:

- 1- Ordonner les échantillons selon l'ordre bit-inversé
- 2- Calculer pour chaque étage, les papillons correspondants.

2.8 FFT pour les signaux réels :

Soit $x(n)$ un signal réel de durée N et $X(k)$ sa TFD. On a :

$$\begin{aligned} X(k) &= \sum_{n=0}^{N-1} x(n) W_N^{nk} \\ &= \sum_{n=0}^{N/2-1} x(2n) \cdot W_{N/2}^{nk} + W_N^k \sum_{n=0}^{N/2-1} x(2n+1) \cdot W_{N/2}^{nk} \end{aligned} \quad (2.23)$$

et soit $z(n)$ un signal complexe tel que :

$$z(n) = x(2n) + j.x(2n+1) \quad , \quad n = 0, \dots, N/2 - 1 \quad (2.24)$$

Sa TFD $Z(k)$ est calculée par :

$$Z(k) = \sum_{n=0}^{N/2-1} [x(2n) + j.x(2n+1)]. W_{N/2}^{nk} \quad (2.25)$$

On obtient :

$$Z\left(\frac{N}{2} - k\right) = \sum_{n=0}^{N/2-1} [x(2n) + j.x(2n+1)]. W_{N/2}^{-nk} \quad (2.26)$$

$$Z^*\left(\frac{N}{2} - k\right) = \sum_{n=0}^{N/2-1} [x(2n) - j.x(2n+1)]. W_{N/2}^{+nk} \quad (2.27)$$

La combinaison de (2.25) et (2.26) donne :

$$\sum_{n=0}^{N/2-1} x(2n). W_{N/2}^{nk} = \frac{1}{2} [Z(k) + Z^*\left(\frac{N}{2} - k\right)] \quad (2.28)$$

$$\text{et} \quad \sum_{n=0}^{N/2-1} x(2n+1). W_{N/2}^{nk} = \frac{1}{2j} [Z(k) - Z^*\left(\frac{N}{2} - k\right)] \quad (2.29)$$

En remplaçant (2.28) et (2.29) dans (2.23), on obtient :

$$X(k) = \frac{1}{2} [Z(k) + Z^*\left(\frac{N}{2} - k\right) - j.W_N^k (Z(k) - Z^*\left(\frac{N}{2} - k\right))] \quad (2.30)$$

Ainsi la FFT d'une séquence réelle de durée N est calculée à partir d'une FFT d'un signal complexe de durée $N/2$. Ceci donne une réduction du nombre d'opérations de près de la moitié.

La FFT du signal réel est calculée comme suit :

1- Constituer le signal complexe selon la relation (2.24), c'est ce qu'on appelle dispersion des données.

2- Calculer la FFT du signal complexe.

3- Reconstituer la TFD du signal réel à partir de celle du signal complexe selon la relation (2.30), c'est ce qu'on appelle rassemblement des données.

2.9 FFT inverse :

La TFD inverse est donnée par la relation (2.13) :

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k). W_N^{-nk}$$

En passant par l'expression conjuguée, on obtient :

$$x^*(n) = \frac{1}{N} \sum_{k=0}^{N-1} X^*(k) \cdot W_N^{+nk} \quad (2.31)$$

La somme a la forme d'une TFD selon la relation (2.12), on obtient donc :

$$x^*(n) = \frac{1}{N} \text{TFD} [X^*(k)]$$

En repassant à l'expression conjuguée, on obtient l'expression de la TFD inverse :

$$x(n) = \frac{1}{N} [\text{TFD}(X^*(k))]^* \quad (2.32)$$

ou, en terme de FFT :

$$x(n) = \frac{1}{N} [\text{FFT}(X^*(k))]^* \quad (2.33)$$

Ceci met en évidence que l'algorithme de la FFT peut être utilisé pour calculer la TFD inverse.

3.1 Introduction :

Un filtre numérique est en général un système linéaire invariant à temps discret réalisé en utilisant une arithmétique à précision finie. Sa conception comporte trois phases :

- 1- Les spécifications désirées du système.
- 2- L'approximation de ces spécifications en utilisant un système discret causal.
- 3- La réalisation du système en utilisant une arithmétique à précision finie.

Les spécifications du système dépendent de l'application considérée, elles se situent en général dans le domaine fréquentiel, comme c'est le cas des filtres sélectifs de fréquence: passe-bas, passe-haut, passe-bande et coupe-bande. Pour ces types de filtres, les spécifications prennent la forme de tolérances [7].

Dans le cas du filtre passe-bas, ces spécifications sont: (voir figure 3.1)

- Une bande passante dans laquelle la réponse fréquentielle doit approximer 1 avec une erreur de $\mp \delta_p$

$$1 - \delta_p \leq |H(f)| \leq 1 + \delta_p, \quad |f| \leq f_p$$

- Une bande coupée dans laquelle la réponse fréquentielle doit approximer 0 avec une erreur δ_s

$$H(f) \leq \delta_s, \quad f_s \leq |f| \leq F/2$$

F étant la fréquence d'échantillonnage.

- Une bande de transition de largeur non nulle ($f_s - f_p$) dans laquelle la réponse fréquentielle décroît de la bande passante à la bande coupée.

La figure 3.2 montre les tolérances pour les autres types de filtres.

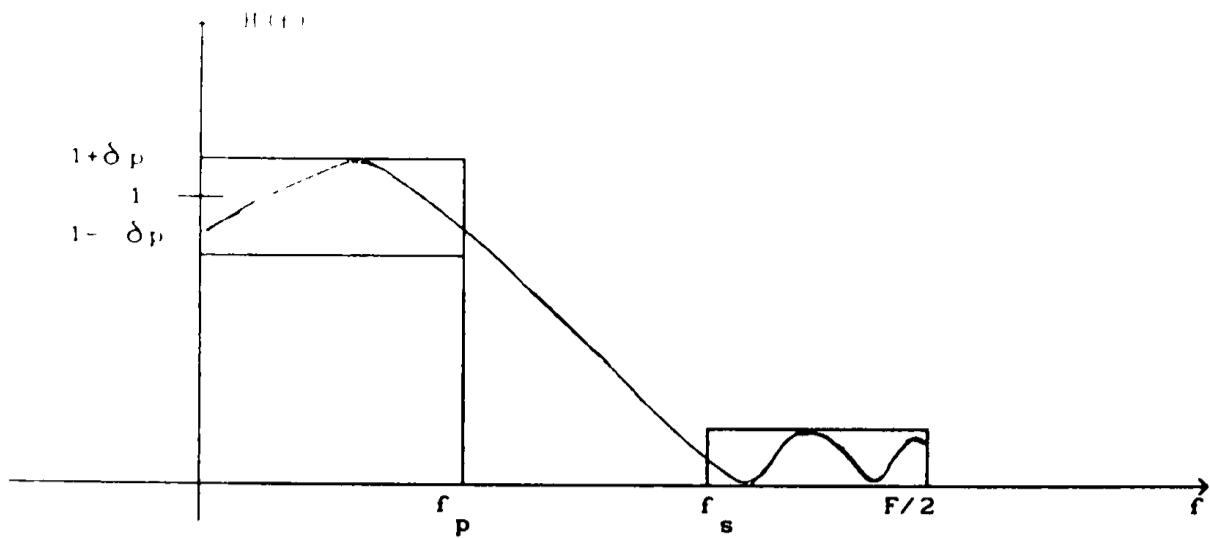


Figure 3.1 : Limite de tolérances d'un filtre passe-bas

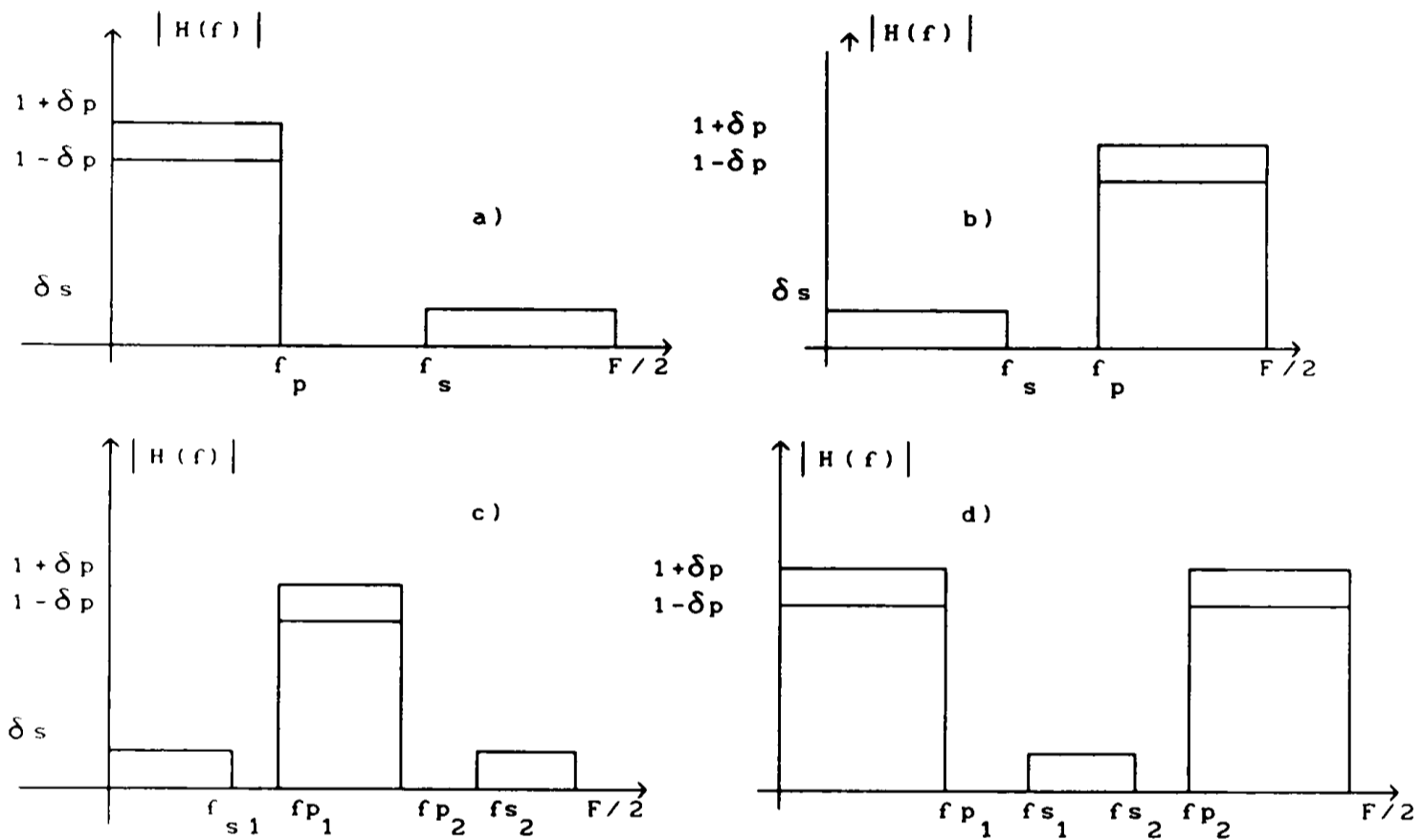


Figure 3.2 : limites de tolérance pour les filtres :

- a) Passe-bas c) Passe-bande
 b) Passe-haut d) Coupe bande

Le problème de l'approximation consiste à chercher un système linéaire invariant dont la réponse fréquentielle respecte les tolérances imposées.

Ce problème fait surgir deux grandes classes de filtres :

- Les filtres récurrents, ou les filtres à réponse impulsionnelle finie (filtres RII ou en anglais IIR). Pour ces filtres l'approximation est faite sous forme d'une fraction rationnelle.

- Les filtres non récurrents, ou filtres à réponse impulsionnelle finie, (filtre RIF ou en anglais FIR). L'approximation est faite sous forme d'un polynôme.

3.2 Filtres IIR :

La fonction de transfert d'un filtre récursif est donnée par:

$$H(z) = \frac{\sum_{i=0}^N A_i \cdot z^{-i}}{1 + \sum_{i=1}^N B_i \cdot z^{-i}} \quad (3.1)$$

La réponse fréquentielle du filtre est obtenue en posant $z = e^{j\omega T}$

$$H(e^{j\omega T}) = \frac{\sum_{i=0}^N A_i \cdot e^{-j\omega T i}}{1 + \sum_{i=1}^N B_i \cdot e^{-j\omega T i}} \quad (3.2)$$

Pour que cette réponse respecte les spécifications désirées, une méthode d'approximation devra être adoptée.

La méthode la plus utilisée pour l'approximation des filtres IIR consiste à transposer les résultats connus de l'approximation des filtres analogiques. C'est la méthode qu'on donnera ici.

D'autres méthodes existent et elles sont dans la majorité basées sur la programmation linéaire. Parmi ces méthodes, on a :

- Minimisation de l'erreur quadratique moyenne [13][17],
- Minimisation de l'erreur L_p [18],
- Placement des pôles et des zéros [6].

3.2.1 Approximations des filtres analogiques :

On donnera une méthode classique d'approximation des filtres passe-bas analogiques, c'est le cas des filtres de Butterworth. On établira les propriétés du filtre passe-bas prototype normalisé dont la pulsation de coupure est égale à l'unité. Les autres filtres en découlent par une transformation appropriée.

3.2.1.1 Les filtres de Butterworth :

La courbe d'amplitude des filtres de Butterworth varie d'une façon monotone dans la bande passante et dans la bande coupée. Le carré de l'amplitude d'un filtre de Butterworth d'ordre N est de la forme :

$$|H(\omega)|^2 = \frac{1}{1 + \epsilon^2 \left(\frac{\omega}{\omega_p}\right)^{2N}} \quad (3.3)$$

La figure (3.3) montre la forme de cette courbe .

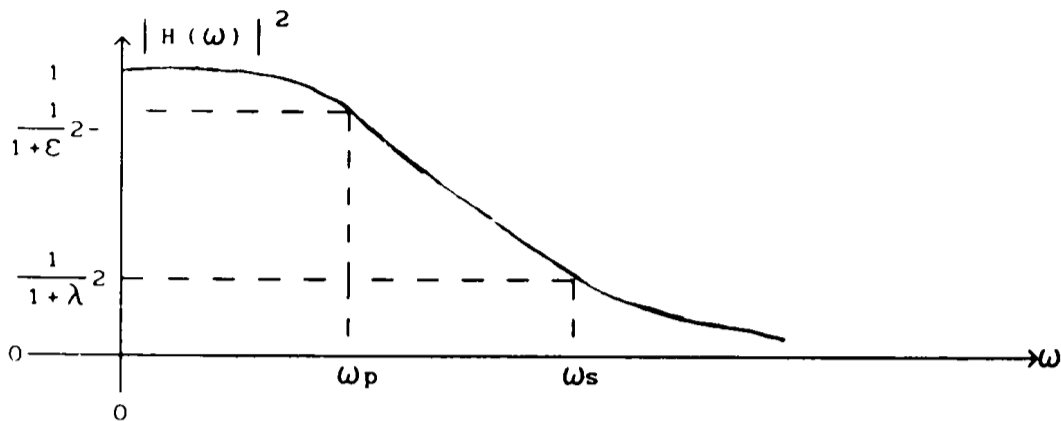


Figure 3.3 : Courbe d'amplitude des filtres de Butterworth

De cette courbe, on définit les paramètres de conception suivants

. Tolérance de la bande passante : $|\omega| \leq \omega_p$

$$|H(\omega)|^2 > \frac{1}{1 + \epsilon^2} \quad (3.4)$$

. Tolérance de la bande coupante : $|\omega| > \omega_s$

$$|H(\omega)| < \frac{1}{1 + \lambda^2} \quad (3.5)$$

En posant $\omega = \omega_s$, et en substituant l'équation (3.3) dans (3.5), on obtient l'expression :

$$\epsilon = (10^{0.1AP} - 1)^{1/2} \quad (3.6.a)$$

$$\lambda = (10^{0.1As} - 1)^{1/2} \quad (3.6.b)$$

$$A = \frac{\lambda}{\epsilon} = \left(\frac{10^{0.1As} - 1}{10^{0.1Ap} - 1} \right)^{0.5} \quad (3.7)$$

et
$$K_0 = \frac{\omega_p}{\omega_s} \quad (3.8)$$

On obtient :

$$N \geq \frac{\log A}{\log \frac{1}{K_0}} \quad (3.9)$$

Les poles normalisés du plan S sont trouvés en posant le dénominateur de l'équation (3.1) égal à zéro. Pour normaliser le résultat, on pose $\omega_p = 1$ et $\epsilon = 1$, alors :

$$1 + \omega^{2N} = 0$$

En posant $S = j\omega$, on obtient :

$$(-S^2)^N + 1 = 0$$

Et en exprimant -1 en notation polaire :

$$(-1)^N S^{2N} = e^{j(2K-1)\pi} = -1 \quad K = 1, 2, \dots, N$$

La $K^{\text{ième}}$ racine pour les pôles dans le demi plan gauche peut être exprimée par :

$$S_k = \sigma_k + j\omega_k = e^{j(2K+N-1)\pi/2N} = je^{j(2K-1)\pi/2N}$$

Finalement, les pôles du filtre passe-bas de Butterworth normalisé peuvent être obtenu par :

$$S_k = -\sin\left(\frac{2K-1}{2N}\pi\right) + j\cos\left(\frac{2K-1}{2N}\pi\right)$$

$$K = \begin{cases} 1, 2, \dots, \frac{N+1}{2} & \text{pour } N \text{ impair} \\ 1, 2, \dots, \frac{N}{2} & \text{pour } N \text{ pair} \end{cases} \quad (3.10)$$

Les équations (3.6 et 3.10) permettent d'obtenir le filtre de Butterworth normalisé à partir des spécifications. Pour obtenir le filtre dénormalisé, on utilise une transformation appropriée. Voir table 3.1 et 3.2.

3.2.1.2 Transformation des bandes de fréquence :

Dans la section précédente, on a traité le problème de l'approximation dans le cas du filtre de Butterworth passe-bas normalisé ($\omega_p = 1$). La table 3.1 résume le passage du filtre passe-bas aux autres types de filtres.

Types de filtre	Transformation
Passe-bas	$S \longrightarrow S/\omega_p$
Passe-haut	$S \longrightarrow \omega_p / S$
Passe-bande	$S \longrightarrow \frac{S^2 + \omega_{p1}\omega_{p2}}{S(\omega_{p2} - \omega_{p1})}$
Coupe bande	$S \longrightarrow \frac{S(\omega_{p2} - \omega_{p1})}{S^2 + \omega_{p2}\omega_{p1}}$

Table 3.1 : Transformation des bandes de fréquence

3.2.2 Transformation bilinéaire :

Les approximations des filtres récursifs sont obtenus à partir de celles des filtres analogiques en utilisant plusieurs méthodes. Parmi celles-ci, on a :

- Méthode de l'invariance impulsionnelle [19],
- Méthode de l'invariance indicielle [19][20],
- Méthode basée sur la solution numérique de l'équation différentielle [7],
- Transformation bilinéaire [14][19][20][21],

Cette dernière méthode est la plus utilisée et c'est celle qu'on présentera .

La transformation bilinéaire est définie par :

$$S = \frac{2}{T} \frac{Z-1}{Z+1} \quad (3.11)$$

Cette transformation permet d'écrire :

$$H_0(z) = H_A(s) \Big|_{s=\frac{2}{T} \frac{z-1}{z+1}} \quad (3.12)$$

où T est la période d'échantillonnage.

Ce qui permet d'obtenir la fonction de transfert du système discret à partir de la fonction de transfert du système analogique. On peut tirer les propriétés de cette transformation comme suit [14]:

avec $s = \sigma + j\omega$

et $z = re^{j\vartheta}$

On obtient :

$$r = \left[\frac{\left(\frac{2}{T} + \sigma\right)^2 + \omega^2}{\left(\frac{2}{T} - \sigma\right)^2 + \omega^2} \right]^{1/2} \quad \text{a)} \quad (3.13)$$

$$\text{et } \vartheta = \tan^{-1}\left(\frac{\omega}{\frac{2}{T} + \sigma}\right) + \tan^{-1}\left(\frac{\omega}{\frac{2}{T} - \sigma}\right) \quad \text{b)}$$

Il est clair que :

- pour $\sigma > 0$, on a $r > 1$
- pour $\sigma = 0$, on a $r = 1$
- pour $\sigma < 0$, on a $r < 1$

Donc on a les propriétés suivantes :

1- Au demi plan gauche du plan s correspond l'intérieur du cercle unité , $|z| = 1$, du plan z .

2- A l'axe imaginaire $j\omega$, du plan s correspond le cercle unité $|z| = 1$, du plan z .

3- Au demi plan droit du plan s correspond l'extérieur du cercle unité, $|z| = 1$, du plan z .

Pour $\vartheta = 0$, $r = 1$, on tire de l'équation (3.13.b) :

$$\vartheta = 2 \tan^{-1}\left(\frac{\omega T}{2}\right) \quad (3.14)$$

Ceci donne :

$$\begin{aligned} \vartheta &= 0 \quad \text{pour } \omega = 0 \\ \vartheta &= \pi \quad \text{pour } \omega \longrightarrow \infty \\ \vartheta &= -\pi \quad \text{pour } \omega \longrightarrow -\infty \end{aligned}$$

L'origine du plan s a pour image le point $(+1,0)$ du plan z . La partie positive et la partie négative du plan s ont pour images respectives le demi cercle supérieur et le demi cercle inférieur du plan z . La transformation est illustrée par la figure (3.4).

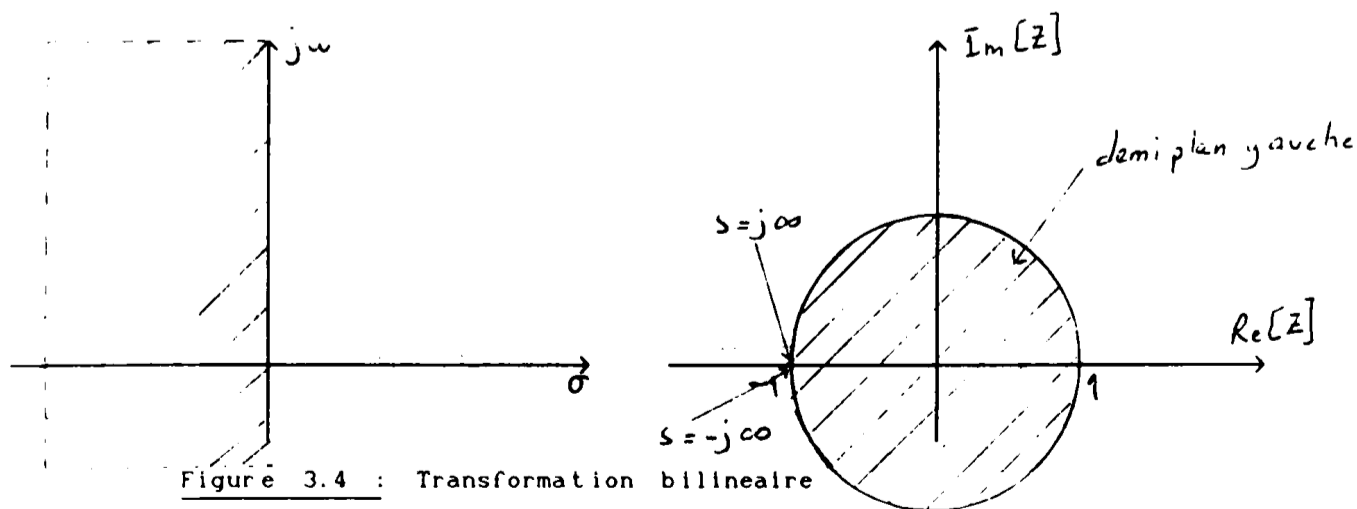


Figure 3.4 : Transformation bilinéaire

Cette transformation a pour effet que le système analogique et le système numérique possèdent le même comportement fréquentiel.

$$\text{Si} \quad M_1 \leq H_A(j\omega) \leq M_2 \quad \text{pour} \quad \omega_1 \leq \omega \leq \omega_2$$

$$\text{alors} \quad M_1 \leq H_D(e^{j\vartheta}) \leq M_2 \quad \text{pour} \quad \vartheta_1 \leq \vartheta \leq \vartheta_2$$

ϑ_1 et ϑ_2 sont liées aux fréquences ω_1 et ω_2 par la relation (3.14). Donc les bandes passantes et coupantes du filtre analogique sont transformées en bandes passantes et coupantes pour le filtre numérique.

3.2.3 Plan de conception

De ce qui précède, on déduit que le plan de conception des filtres IIR est le suivant :

- 1- On élabore les spécifications du filtre.
- 2- On déduit les paramètres du filtre passe-bas prototype pour l'approximation choisie en utilisant la table 3.2
- 3- On tire la fonction de transfert du filtre passe-bas analogique prototype.
- 4- On tire la fonction de transfert du filtre analogique correspondant en utilisant la table 3.1.
- 5- On tire la fonction de transfert du filtre numérique en utilisant la transformation bilinéaire - équation 3.11

Approximation	Equation de l'ordre du Filtre
Butterworth	$N > \frac{\text{Log} A}{\text{Log}(1/k)}$
Passe-bas	Passe-haut
$k = k_0 = \frac{\tan(\pi f_p / F)}{\tan(\pi f_s / F)}$	$k = \frac{1}{k_0}$
Passe-bande	Coupe-bande
$k = \begin{cases} k_1, & \text{si } k_c \geq k_B \\ k_2, & \text{si } k_c < k_B \end{cases}$	$k = \begin{cases} 1/k_2 & \text{si } k_c \geq k_B \\ 1/k_1 & \text{si } k_c < k_B \end{cases}$
$k_A = \tan(\pi f_{p_2} / F) - \tan(\pi f_{p_1} / F)$	
$k_B = \tan(\pi f_{p_1} / F) \cdot \tan(\pi f_{p_2} / F)$	
$k_C = \tan(\pi f_{s_1} / F) \cdot \tan(\pi f_{s_2} / F)$	
$k_1 = \frac{k_A \tan(\pi f_{s_1} / F)}{k_B - \tan^2(\pi f_{s_2} / F)}$	$k_2 = \frac{k_A \tan(\pi f_{s_2} / F)}{\tan^2(\pi f_{s_1} / F) - k_B}$

Table 3.2 : Equations donnant l'ordre du filtre numerique

3.3 Filtres FIR :

La fonction de transfert d'un filtre FIR causal est donnée par [6]:

$$H(z) = \sum_{n=0}^{N-1} h(n) \cdot z^{-n} \quad (3.15)$$

où $h(n)$ est la réponse impulsionnelle du filtre. La transformée de Fourier de la séquence $h(n)$ est donnée par :

$$H(e^{j\omega T}) = \sum_{n=0}^{N-1} h(n) \cdot e^{-j\omega T n} = |H(e^{j\omega T})| \cdot e^{j\vartheta(\omega)} \quad (3.16)$$

L'amplitude et la phase sont définies par :

$$M(\omega) = |H(e^{j\omega T})| \quad (3.17)$$

$$\vartheta(\omega) = \tan^{-1} \frac{\text{Im } H(e^{j\omega T})}{\text{Re } H(e^{j\omega T})}$$

On définit le retard de phase et le retard de groupe du filtre comme suit :

$$\tau_p = -\frac{\vartheta(\omega)}{\omega} \quad \text{et} \quad \tau_g = -\frac{d\vartheta(\omega)}{d\omega} \quad (3.18)$$

Les filtres pour lesquels τ_p et τ_g sont constantes sont appelés filtres à retard constant ou filtres à phase linéaire. Ils constituent la classe la plus importante des filtres FIR.

Pour que la phase soit linéaire, on doit avoir :

$$\vartheta(\omega) = -\tau.\omega \quad -\pi < \omega < \pi \quad (3.19)$$

Des équations (3.16), (3.17) et (3.19), la réponse de phase peut être exprimée par :

$$\vartheta(\omega) = -\tau.\omega = \tan^{-1} - \frac{\sum_{n=0}^{N-1} h(n) \cdot \sin(\omega nT)}{\sum_{n=0}^{N-1} h(n) \cdot \cos(\omega nT)} \quad (3.20)$$

où

$$\tan \omega\tau = \frac{\sum_{n=0}^{N-1} h(n) \cdot \sin(\omega nT)}{\sum_{n=0}^{N-1} h(n) \cdot \cos(\omega nT)} \quad (3.21)$$

Finalement, on obtient :

$$\sum_{n=0}^{N-1} h(n) \cdot \sin(\omega\tau - \omega nT) = 0 \quad (3.22)$$

On montre [22] que la solution de (3.22) est donnée par :

$$\tau = \frac{(N-1)T}{2} \quad (3.23)$$

et

$$h(n) = h(N-1-n) \quad \text{pour } 0 \leq n \leq (N-1) \quad (3.24)$$

De l'équation (3.24), on voit que la réponse impulsionnelle est symétrique selon que N est pair ou impair, le centre de symétrie survient entre deux échantillons ou coïncide avec un échantillon, ceci est illustré à la figure 3.5

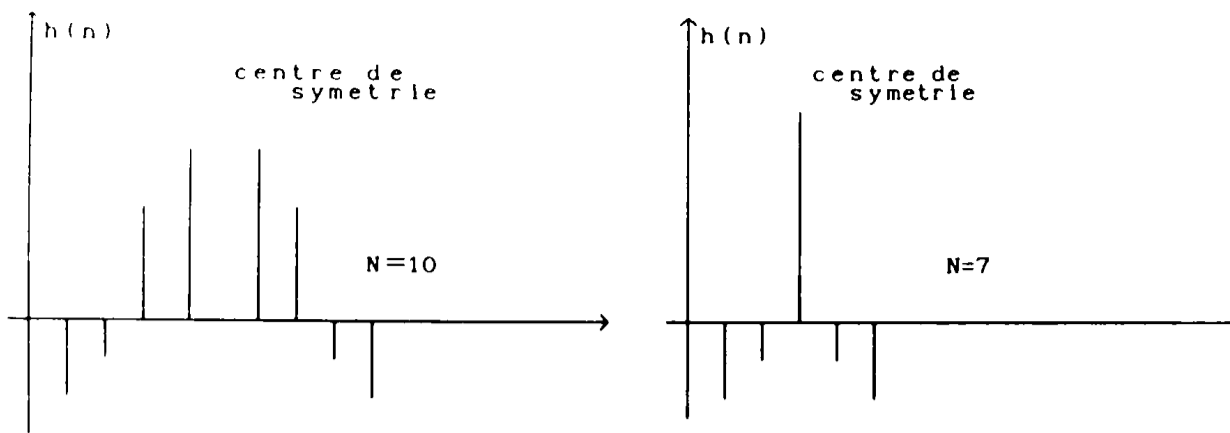


Figure 3.5 : Réponse impulsionnelle d'un filtre à phase linéaire

3.3.1 Réponse fréquentielle des filtres FIR à phase linéaire

La réponse fréquentielle des filtres FIR à phase linéaire est obtenue en exploitant la relation 3.24 dans l'équation (3.16). Deux cas sont à considérer :

1- Pour N impair, on obtient :

$$H(e^{j\omega T}) = e^{-j\omega(N-1)T/2} \left\{ \sum_{n=0}^{\frac{N-1}{2}} a(n) \cdot \cos(\omega n T) \right\} \quad (3.25)$$

$$\text{avec } \begin{cases} a(0) = h\left(\frac{N-1}{2}\right) \\ a(n) = 2 h\left(\frac{N-1}{2} - n\right) \end{cases}, \quad n = 1 \dots \frac{N-1}{2}$$

2- Pour N pair, on obtient :

$$H(e^{j\omega T}) = e^{-j\omega(N-1)T/2} \left\{ \sum_{n=0}^{N/2} b(n) \cdot \cos\left(\omega \left(\frac{n-1}{2}\right)T\right) \right\} \quad (3.26)$$

$$\text{avec } b(n) = 2 h\left(\frac{N}{2} - n\right), \quad n = 1 \dots \frac{N}{2}$$

3.3.2 Calcul des filtres FIR par développement en série de Fourier:

On sait d'après le premier chapitre que la représentation fréquentielle d'un signal discret est périodique. Cette représentation est donc décomposable en série de Fourier.

Donc, la réponse fréquentielle désirée du filtre peut être exprimée par [6]:

$$H(e^{j2\pi fT}) = \sum_{n=-\infty}^{+\infty} h_d(n) \cdot e^{j2\pi n f T} \quad (3.27)$$

où les coefficients de Fourier $h_d(n)$ sont la réponse impulsionnelle désirée du filtre. Ils sont calculés par :

$$h_d(n) = \frac{1}{N} \int_{-F/2}^{F/2} H(e^{j2\pi fT}) \cdot e^{j2\pi n f T} df \quad (3.28)$$

En posant $z = e^{j2\pi fT}$ dans l'équation (3.27), on obtient la fonction de transfert du filtre :

$$H(z) = \sum_{n=-\infty}^{+\infty} h_d(n) \cdot z^{-n} \quad (3.29)$$

On voit que les coefficients $h_d(n)$ sont de durée infinie et la fonction de transfert représente un système non causal. Donc, on doit limiter la durée des coefficients en les multipliant par une fenêtre de troncature, et rendre le système causal en multipliant $H(z)$ par $z^{-(N-1)/2}$, ce qui donne :

$$H'(z) = z^{-(N-1)/2} \cdot \sum_{n=-\frac{N-1}{2}}^{\frac{N-1}{2}} h_d(n) \cdot \omega(n) \cdot z^{-n} = z^{-\frac{N-1}{2}} \cdot H(z) \quad (3.30)$$

où $\omega(n)$ représente la fenêtre utilisée.

La multiplication par $z^{-(N-1)/2}$ n'a aucun effet sur la courbe d'amplitude du filtre. Le retard de groupe augmentera d'une constante $(\frac{N-1}{2})T$

Cette méthode est la plus simple et la plus utilisée.

D'autres méthodes basées sur la solution numérique du problème sont utilisées [23-26]. Elles aboutissent à la solution optimale après plusieurs itérations.

3.3.3 Caractéristiques des filtres FIR :

Certains des avantages des filtres FIR par rapport aux filtres IIR sont comme suit :

1- Les filtres FIR peuvent être conçus exactement avec une phase linéaire. La phase linéaire est importante pour les applications

où les distorsions de phase dues à la non linéarité peuvent dégrader les performances -par exemple, pour le traitement de la parole ou pour la transmission des données.

2- Les filtres FIR sont toujours stables vu que leur fonction de transfert ne possède pas de pôles sauf en $z=0$.

3- Le bruit de quantification dû à l'arithmétique à précision finie peut être rendu négligeable.

Parmi les désavantages, on a:

1- L'ordre d'un filtre FIR est beaucoup plus grand que celui d'un filtre IIR qui réalise la même courbe de fréquence.

2- Il n'existe pas d'équations générales pour la conception des filtres FIR. Plusieurs des méthodes utilisées sont itératives et exigent un outil de calcul puissant pour leur implantation.

3.4 Structures des filtres numériques :

3.4.1 Filtres IIR :

La fonction de transfert d'un filtre IIR est une fraction rationnelle de la forme [11]:

$$H(z) = \frac{\sum_{l=0}^N a_l \cdot z^{-l}}{1 + \sum_{l=0}^N b_l \cdot z^{-l}} \quad (3.31)$$

à laquelle correspond la structure directe canonique de la figure (3.6).

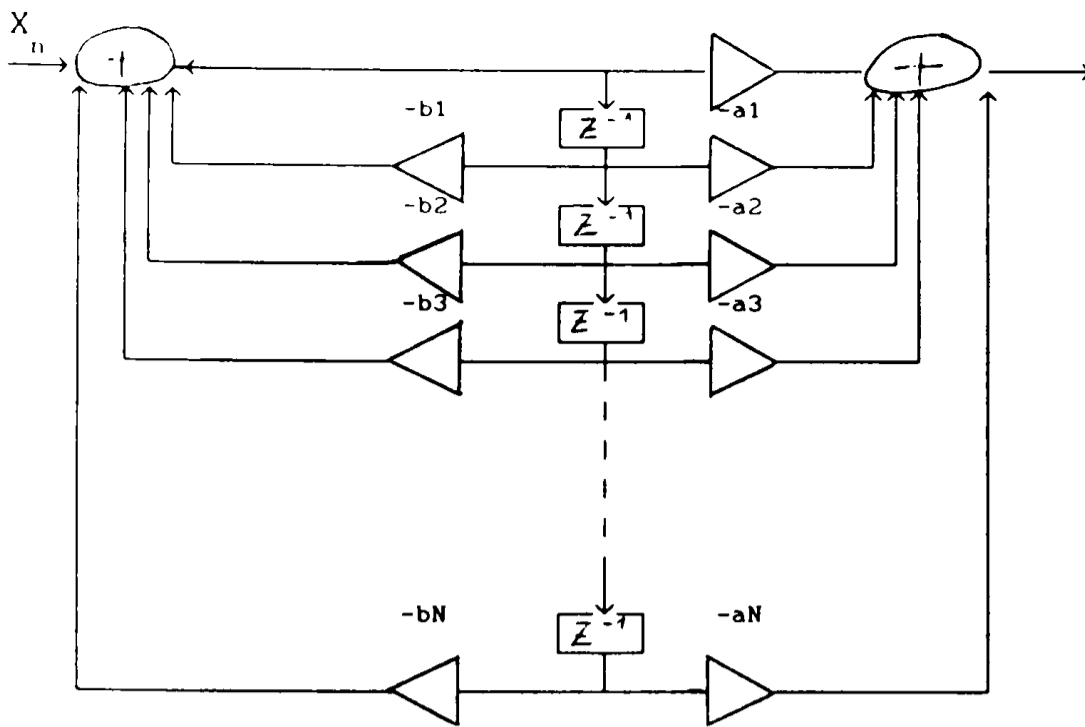


Figure 3.6 : Structure Directe Canonique

Si on factorise $H(z)$ sous la forme [11]:

$$H(z) = k \prod_{i=1}^q \frac{1 + a_{1i} z^{-1} + a_{2i} z^{-2}}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \quad (3.32)$$

on obtient la forme canonique en cascade, de la figure (3.7)

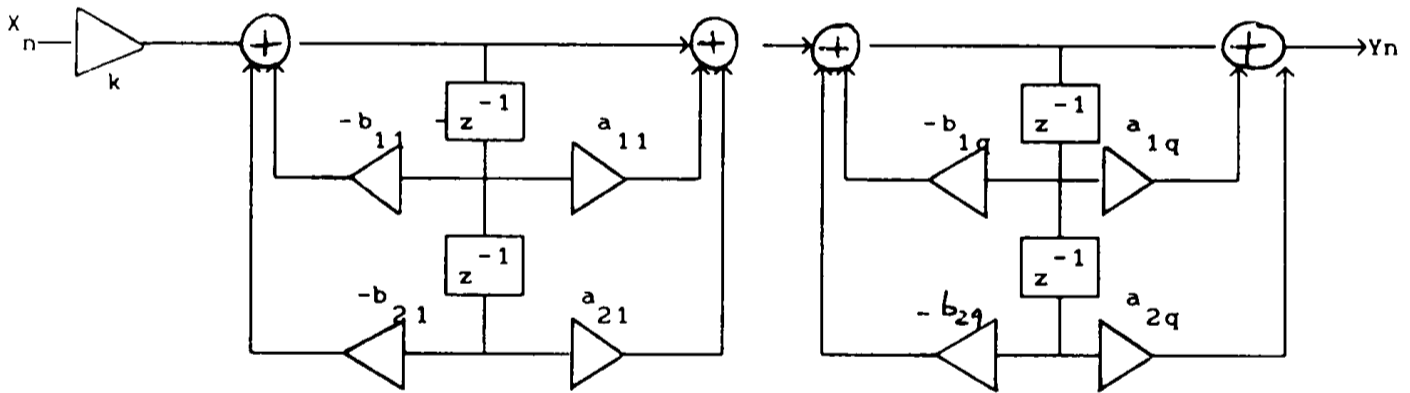


Figure 3.7 : Structure canonique en cascade

Tandis que la structure canonique en parallèle correspond à la décomposition en fraction simples (fig 3.10) [11]:

$$H(z) = k + \sum_{i=1}^q \frac{h_{0i} + h_{1i} z^{-1}}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \quad (3.33)$$

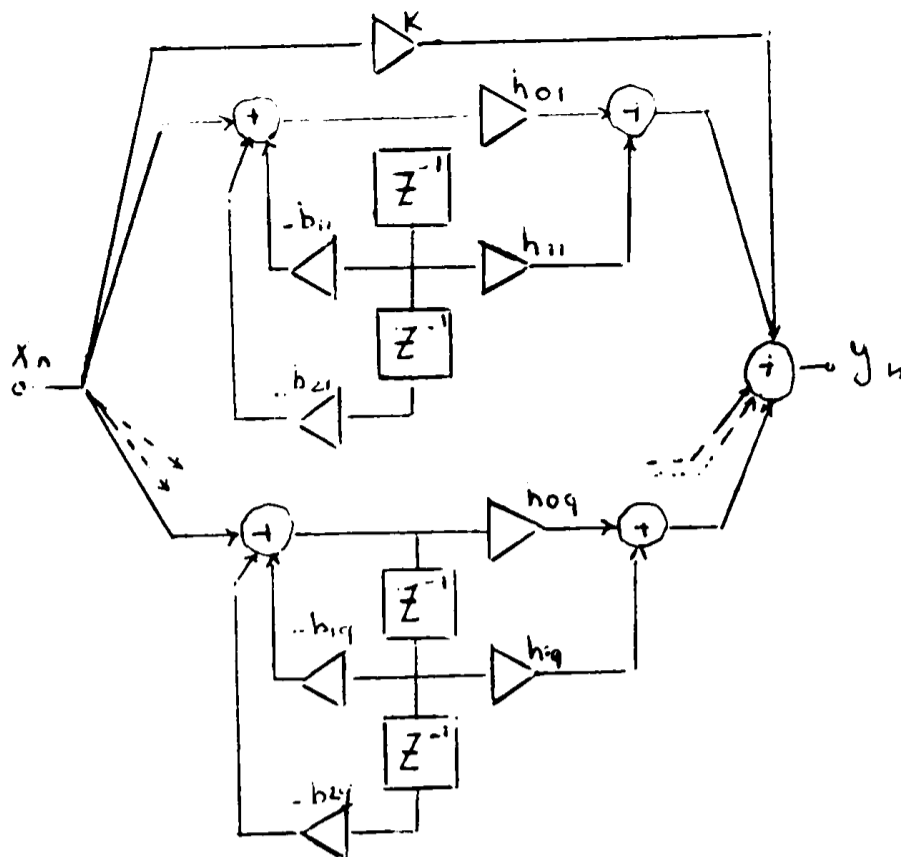


Figure 3.8 : Structure canonique en parallele

3.4.1.1. Sensibilité des filtres IIR.

Quelle que soit la réalisation physique d'un filtre numérique, les paramètres qui définissent sa fonction de transfert sont représentés avec une précision finie. Le choix d'une structure particulière dépend donc de la sensibilité de la fonction de transfert aux erreurs de quantification commises sur les coefficients [11]

Soient $H(z)$ la fonction de transfert d'un filtre numérique, $\{X\}$ le vecteur des coefficients, $\{\Delta X\}$ celui des écarts séparant les coefficients quantifiés de leur valeur idéale et $A(\omega)$ l'affaiblissement défini par [11] :

$$A(\omega) = -20 \text{ Log } |H(z)| \Big|_{z=e^{j\omega}} \quad (3.34)$$

Les écarts $\{\Delta X\}$ provoquent une variation $\Delta A(\omega)$ de l'affaiblissement $A(\omega)$. Cette variation vaut :

$$\Delta A(\omega) = \sum_k \frac{\partial A(\omega)}{\partial X_k} \quad (3.35)$$

On définit la dérivée partielle :

$$S_{X_k}(\omega) = \frac{\partial A(\omega)}{\partial X_k} \quad (3.36)$$

comme étant la sensibilité de l'affaiblissement par rapport au coefficient X_k . De (3.34), il vient [11]:

$$S_{X_k}(\omega) = -8,6859 \operatorname{Re} \left\{ \frac{1}{H(z)} \cdot \frac{\partial H(z)}{\partial X_k} \Big|_{z=e^{j\omega}} \right\} \quad (3.37)$$

La constante -8,6859 correspond à $(-20/\ln(10))$

L'équation (3.35) devient :

$$\Delta A(\omega) = \sum_k S_{X_k}(\omega) \cdot \Delta X_k \quad (3.38)$$

Dans la représentation à virgule fixe, on a :

$$|\Delta X_k| \leq \frac{1}{2} 2^{-b} \quad (3.39)$$

où b désigne le nombre de bits réservés à la partie fractionnaire des coefficients.

Pour pouvoir comparer les différentes structures réalisant la même fonction de transfert, on utilise une méthode statistique. On fait les hypothèses suivantes [11] :

1- L'erreur d'arrondi sur un coefficient X_k est une variable aléatoire uniformément répartie entre $\frac{-q}{2}$ et $\frac{+q}{2}$. Sa variance vaut

$$\sigma_{X_k}^2 = \frac{q^2}{12}, \quad q = 2^{-b} \quad (3.40)$$

2- Les diverses erreurs ΔX_k sont statistiquement indépendantes. La variance de la perturbation $\Delta A(\omega)$ peut alors s'écrire [11]

$$\sigma_{\Delta A}^2 = \sum_k S_{X_k}^2 \cdot \sigma_{X_k}^2 \quad (3.41)$$

et l'écart quadratique moyen vaut :

$$\begin{aligned}\sigma_{\Delta_A} &= P(\omega) \cdot q / \sqrt{12} \\ &= P(\omega) \cdot 2^{-b} / \sqrt{12}\end{aligned}\quad (3.42)$$

La fonction $P(\omega)$ est appelée indice de sensibilité quadratique. C'est sur cette fonction que se base la comparaison des structures.

On peut tirer les sensibilités des structures directe et cascade pour pouvoir les comparer.

Soient b_n ($n = 1 \dots N$) un des coefficients du dénominateur de la forme directe et b_{1k} ($l = 1, 2 ; k = 1 \dots \frac{N}{2}$) un coefficient du dénominateur de la forme cascade. On tire les sensibilités [11] :

$$S_{b_n} = 8,6859 \cdot \text{Re} \left\{ \frac{z^{-n}}{1 + \sum_{l=1}^N b_l z^{-l}} \Big|_{z=e^{j\omega}} \right\} \quad (3.42)$$

$$S_{b_{1k}} = 8,6859 \cdot \text{Re} \left[\frac{z^{-1}}{1 + b_{1k} z^{-1} + b_{2k} z^{-2}} \Big|_{z=e^{j\omega}} \right] \quad (3.43)$$

En se basant sur la figure 3.9 et en considérant que les Z_{pi} ($i=1 \dots N$) sont les pôles de la fonction de transfert. On peut écrire (4.42) et (4.43) sous les formes :

$$S_{b_n} = 8,6859 \cdot \text{Re} \left\{ \frac{1}{\prod_{l=1}^N |e^{j\omega} - Z_{pl}|} \cdot \text{Cos}(\sum_{l=1}^N \psi_l + \phi) \right\} \quad (3.44)$$

$$S_{b_{1k}} = 8,6859 \cdot \text{Re} \left\{ \frac{1}{|e^{j\omega} - Z_{p1}| |e^{j\omega} - Z_{pk}^*|} \cdot \text{Cos}(\psi_k + \psi_{k^*} + \phi) \right\} \quad (3.45)$$

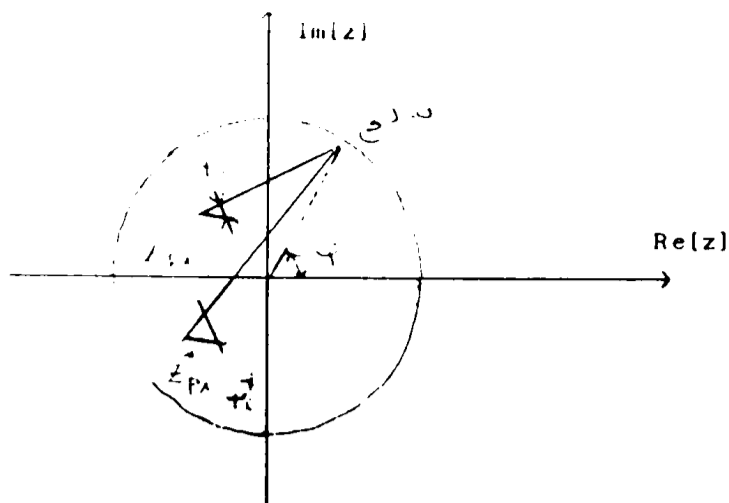


Figure 3.9 : Position de deux poles conjugués dans le plan z

La quantité $|e^{j\omega} - z_{p1}|$ est en général plus petite que l'unité, ceci rend clair que :

$$S_{b_n}^{\max} \gg S_{b_{lk}}^{\max} \quad (3.46)$$

et ce d'autant plus que le degré du filtre est plus élevé [11]:

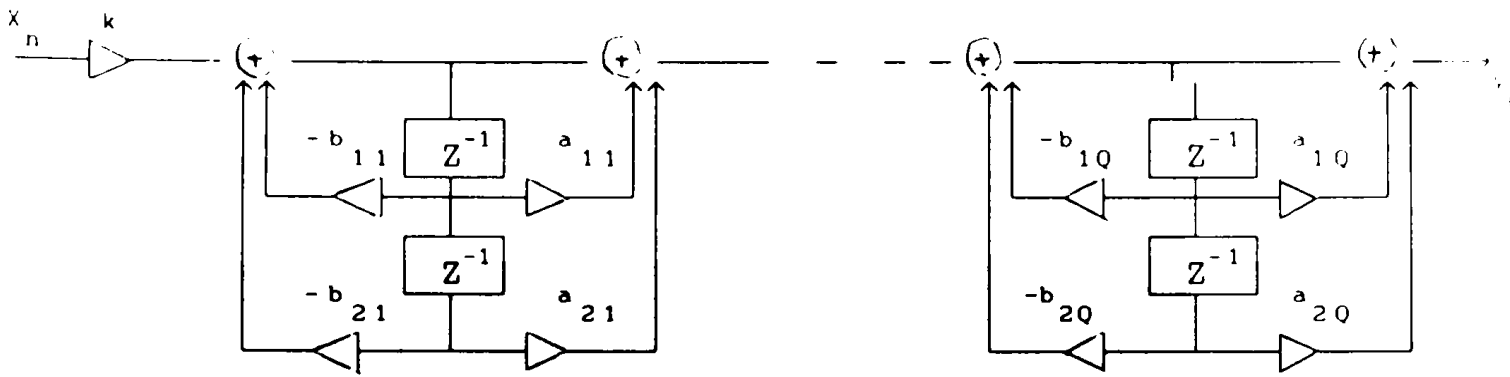
Des résultats semblables peuvent être obtenus pour les coefficients du numérateur.

Donc la structure cascade est beaucoup moins sensible aux erreurs commises sur les coefficients que la structure directe. Ceci explique pourquoi la structure cascade est la plus utilisée. Des résultats pareils à (3.46) peuvent être obtenus pour la structure parallèle.

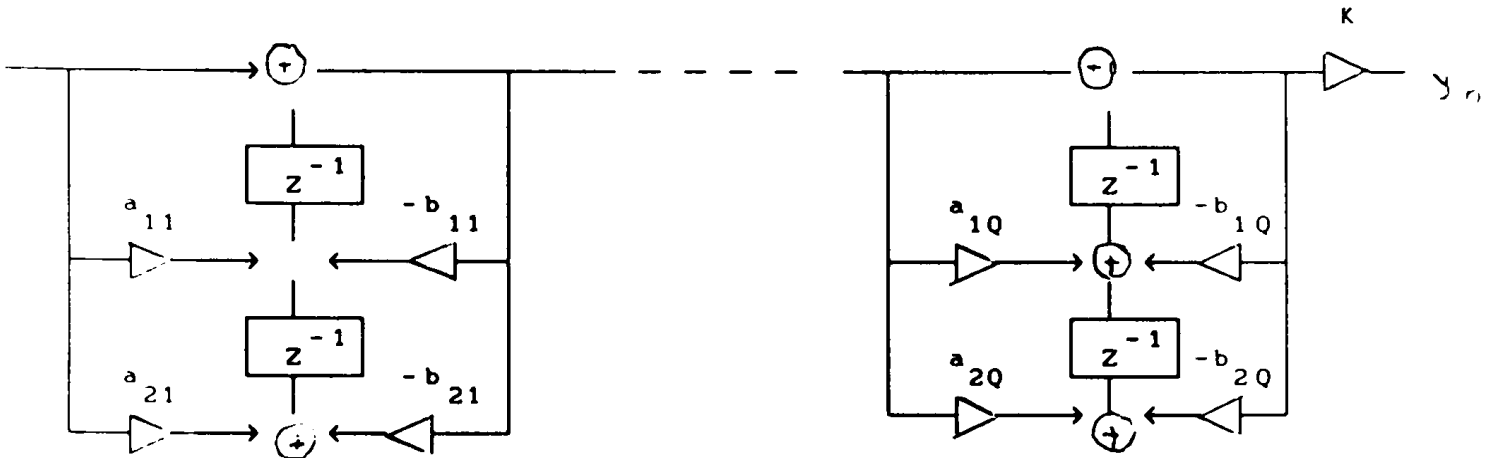
La structure cascade correspond à la forme factorisée de la fonction de transfert :

$$H(z) = K \cdot \prod_{i=1}^N \frac{1 + a_{1i} z^{-1} + a_{2i} z^{-2}}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \quad (3.47)$$

Elle correspond à la mise en cascade de plusieurs cellules du second ordre. Ces dernières peuvent être organisées de plusieurs manières mais les formes les plus utilisées sont les structures 1D et 2D [34] qui correspondent à la figure (3.10).



a) Forme 1D



b) Forme 2D

Figure 3.10 : Structure cascade

3.4.2 Filtres FIR :

La structure la plus utilisé pour les filtres FIR est la structure directe. Elle correspond à la figure (3.11).

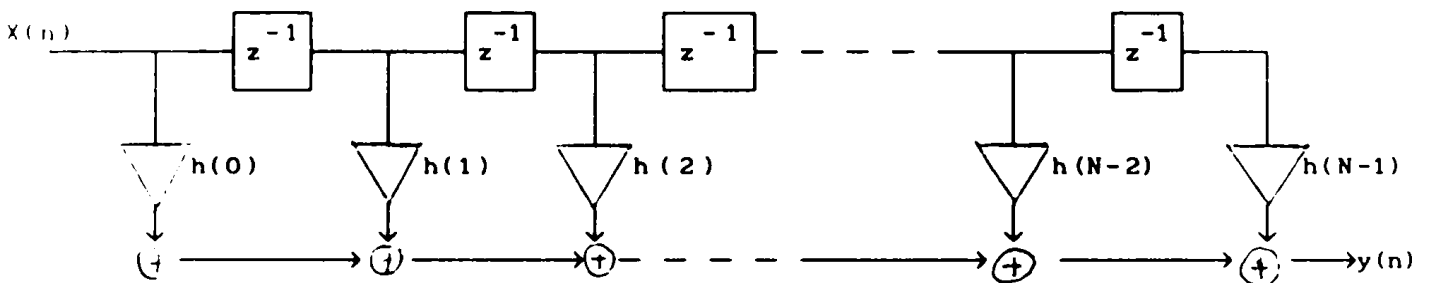


Figure 3.11 : Structure Directe d'un filtre FIR

3.4.2.1 Sensibilité des filtres FIR :

L'amplitude $A(\omega)$ de la réponse fréquentielle d'un filtre FIR est obtenue à partir de l'équation (3.25) par [11]:

$$A(\omega) = \sum_{n=0}^{(N-1)/2} 2h_n \cos\left(\frac{N-1}{2} - n\right)\omega + h\left(\frac{N-1}{2}\right) \quad (3.47)$$

où les $h(n)$ sont les coefficients du filtre et N son ordre. On considère que N est impair .

La sensibilité de l'amplitude par rapport aux coefficients vaut [6]:

$$S_{h_n} = \frac{\partial A(\omega)}{\partial h_n} = 2 \cos\left(\frac{N-1}{2} - n\right)\omega, \quad n = 1 \dots \frac{N-3}{2}$$

$$S_{h_{\frac{N-1}{2}}} = 1 \tag{3.48}$$

On remarque que ces sensibilités sont indépendantes des coefficients du filtre et de plus, elles sont comprises entre -2 et +2.

Une borne supérieure sur la variation $\Delta A(\omega)$ de l'amplitude donnée par :

$$|\Delta A(\omega)| \leq N \cdot \frac{g}{2} \tag{3.49}$$

Ceci explique l'intérêt porté à la structure directe. En plus les filtres FIR réalisés en structure cascade exigent des cellules du quatrième ordre et ne préservent pas la phase linéaire [7].

REGISTRES EN TRAITEMENT

NUMERIQUE DU SIGNAL

4.1 Introduction :

Les algorithmes du traitement numérique du signal, tels que le filtrage numérique et la FFT, sont réalisés soit avec un équipement numérique hardware adéquat, soit par des programmes exécutés sur ordinateur. Dans les deux cas, les coefficients et les données sont stockés sous une forme binaire avec des registres de longueur finie.

Les paramètres des filtres numériques conçus par une des méthodes du troisième chapitre sont obtenus avec une grande précision. Une fois ces paramètres quantifiés, la réponse fréquentielle du filtre résultant sera différente de celle du filtre conçu et peut même dégrader les performances de sorte que la réponse fréquentielle dépasse les tolérances désirées.

La contrainte de la longueur finie des mots (registres) se manifeste de plusieurs façons suivant la représentation choisie: en virgule fixe ou en virgule flottante. Puisque l'ADSP2100 est un processeur à virgule fixe, on ne présentera ici que l'effet de la représentation en virgule fixe, et on commence par définir cette représentation.

4.2 Représentation en virgule fixe :

Un nombre est représenté par une séquence de chiffres binaires (bits) qui sont soit 1, soit 0. Cette séquence est divisée en une partie entière et une partie fractionnaire par un point binaire.

La représentation en virgule fixe est celle pour laquelle la position du point binaire est fixe [7]. Si Δ dénote la position du point, le nombre binaire $(1001_{\Delta}0110)$ a une valeur décimale de 9.375.

On suppose que les registres ont une longueur de $(b+1)$ bits et que le nombre réel x à représenter est tel que [30]:

$$0 \leq |x| \leq 1 \quad (4.1)$$

Si x est tel que $|x| > 1$, on peut normaliser x en décalant le point binaire de L positions tel que l'on ait [30]:

$$x' = 2^L \cdot x \quad (4.2)$$

Le nombre x est représenté par une version quantifiée $Q[x]$ entachée d'une erreur [30]:

$$e = Q[x] - x \quad (4.3)$$

Cette erreur est une variable aléatoire dont les caractéristiques varient d'une représentation à une autre. Il y a trois façons de représenter les nombres en virgules fixe [31]:

- 1- En signe et valeur absolue
- 2- En complément à deux.
- 3- En complément à un.

4.2.1 Représentation en signe et valeur absolue :

Dans cette représentation, les nombres sont exprimés par [30]:

$$Q[x] = (S \ m_1 \ m_2 \ m_3 \ \dots \ m_b)_2 \quad (4.4)$$

où $Q[x]$ est la version quantifiée de x .

S est le bit Signe.

$S = 0$ pour les x positifs

$S = 1$ pour x négatif

Les m_i sont les bits de la magnitude :

$$(\Delta \ m_1 \ m_2 \ \dots \ m_b)_2 = |Q[x]|$$

Il y a deux façons de réaliser la quantification soit par troncature soit par arrondi.

a) Troncature

On prend le nombre x , on le convertit en une fraction binaire et on tronque cette fraction à b bits, comme suit [30]:

$$\begin{aligned} |x| &= (\Delta \ m_1 \ m_2 \ \dots \ m_b \ m_{b+1} \ \dots)_2 \\ |Q[x]| &= (\Delta \ m_1 \ m_2 \ \dots \ m_b)_2 \end{aligned} \quad (4.5)$$

L'indice T dénote la troncature. Pour $x \geq 0$, on a $|x| \geq Q_T[x]$ et l'erreur introduite est [30]:

$$\begin{aligned} E_T &= Q_T[x] - x = |Q_T[x]| - |x| \\ &= - (00 \dots 0 m_{b+1} m_{b+1} \dots)_2 \\ &= - 2^{-b} (m_{b+1} m_{b+1} \dots)_2 \end{aligned} \quad (4.6)$$

où le nombre $(\Delta m_{b+1} m_{b+1} \dots)_2$ est borné par :

$$0 \leq (\Delta m_{b+1} m_{b+2} \dots)_2 < 1 \quad (4.7)$$

Donc pour $x > 0$, on a :

$$0 \geq E_T > - 2^{-b}, \quad x \geq 0 \quad (4.8)$$

Pour $x < 0$, on aura $x \leq Q_T[x]$. On trouve [30]:

$$E_T = - | | Q_T[x]| - |x| | \quad (4.9)$$

et

$$0 \leq E_T < 2^{-b}, \quad x < 0 \quad (4.10)$$

La caractéristique de quantification de $Q_T[x]$ est illustrée par la figure (4.1.a). La densité de probabilité de E_T est montrée à la figure (4.2.a) [30].

La variance du bruit de troncature est donnée par :

$$\sigma_{E_T}^2 = \frac{5}{24} \cdot 2^{-2b} \quad (4.11)$$

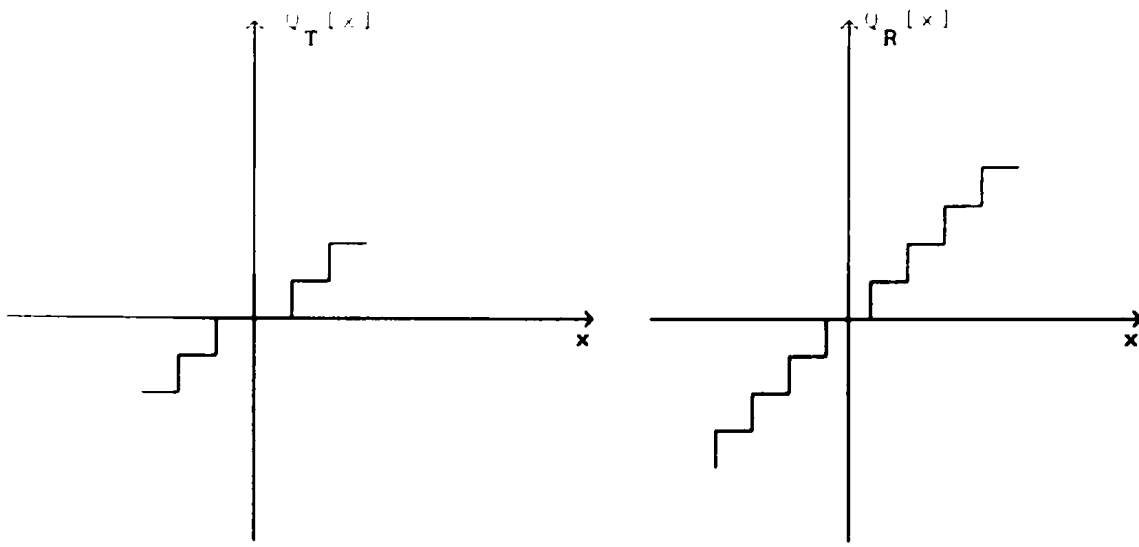


Figure 4.1 : Caractéristique du quantificateur de la représentation
 Signe et valeur absolue
 a) Par troncature b) Par rondi

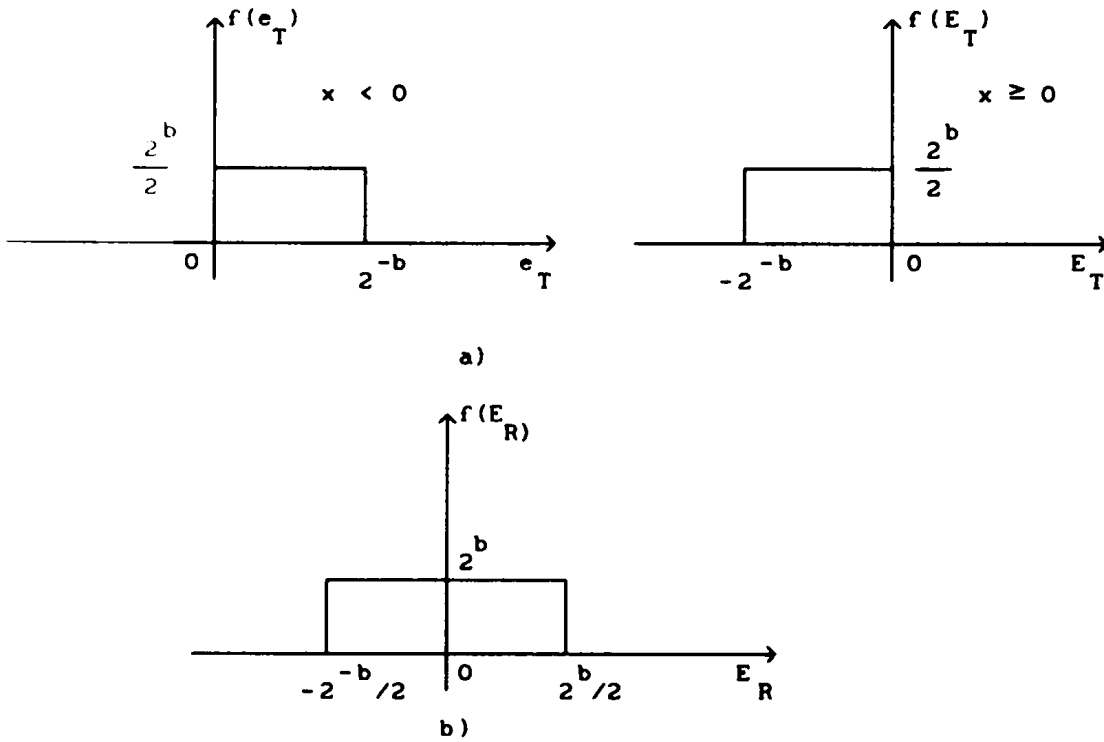


Figure 4.2 : Fonction densité de probabilité de l'erreur de
 quantification de la représentation signe et valeur absolue.
 a) Par troncature b) Par arrondi

b) Arrondi

On prend la valeur absolue de x , à laquelle on ajoute 2^{-b-1} .
 On prend ensuite le résultat et on le tronque à b bits [30].

$$\begin{aligned}
 |x| &= (\Delta \ n_1 \ n_2 \ \dots \ n_b \ n_{b+1} \ \dots)_2 \\
 x &= (S \Delta \ n_1 \ n_2 \ \dots \ n_b \ n_{b+1} \ \dots)_2 \\
 + 2^{-b-1} &= (0 \ \Delta \ 0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots)_2
 \end{aligned}
 \tag{4.12}$$

$$x + 2^{-b-1} = (S \Delta \ m_1 \ m_2 \ \dots \ m_b \ m_{b+1} \ \dots)_2$$

Il en suit la valeur quantifiée de x ,

$$Q_R[x] = (S_{\Delta} n_1 n_2 \dots n_b)_2 \quad (4.13)$$

L'indice R dénote l'arrondi.

L'erreur introduite est :

$$E_R = Q_R[x] - x \quad (4.14)$$

On suppose [30][31][7] que l'erreur introduite est bornée par l'intervalle :

$$\frac{-2^{-b}}{2} \leq E_R \leq \frac{2^{-b}}{2} \quad (4.15)$$

La variance du bruit d'arrondi est [30]:

$$\sigma_{E_R}^2 = \frac{2^{-2b}}{12} \quad (4.16)$$

La caractéristique du quantificateur $Q_R[x]$ est illustrée par la figure (4.1.b). La densité de probabilité de l'erreur introduite est illustrée à la figure (4.2.b).

4.2.2. Représentation en complément à 2 :

Dans cette représentation, les nombres sont exprimés par [30]:

$$\begin{aligned} Q[x] &= (0_{\Delta} m_1 m_2 \dots m_b) & 0 \leq x \leq 1 \\ &= (1_{\Delta} n_1 n_2 \dots n_b) & -1 \leq x \leq 0 \end{aligned} \quad (4.17)$$

où

$$\text{pour } x \geq 0 \quad (\Delta m_1 m_2 \dots m_b) = |Q[x]|$$

$$\begin{aligned} \text{pour } x < 0 \quad (\Delta n_1 n_2 \dots n_b) &= \text{Complément à 2 de } |Q[x]| \\ &= 1 - |Q[x]| \end{aligned}$$

a) Troncature

Pour les nombres positifs, le résultat est identique à la

représentation en signe et valeur absolue [30]:

$$-2^{-b} < E_T < 0 \quad , \quad x \geq 0 \quad (4.18)$$

Pour les nombres négatifs, on a :

$$-2^{-b} < E_T \leq 0 \quad , \quad x < 0 \quad (4.19)$$

Donc pour tout x, on a :

$$-2^{-b} < E_T \leq 0 \quad (4.20)$$

La variance du bruit de quantification est [30]:

$$\sigma_{E_T}^2 = \frac{2^{-b}}{12} \quad (4.21)$$

La caractéristique de ce quantificateur est illustrée par la figure (4.3.a). La densité de probabilité de l'erreur E_T est illustrée à la figure (4.4.a).

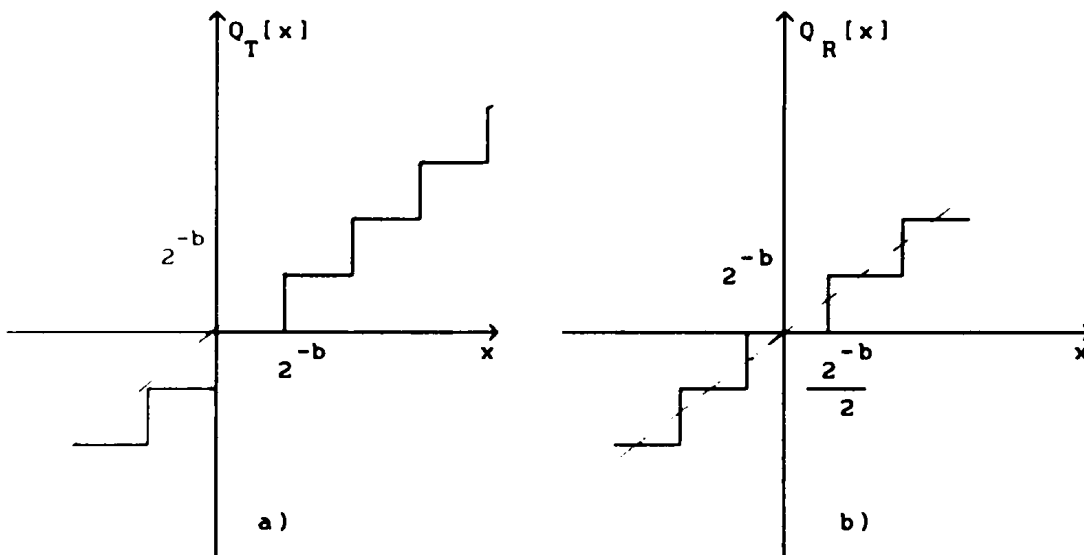


Figure 4.3 : Caractéristique du quantificateur de la représentation en complément à 2
a) Par troncature
b) Par arrondi

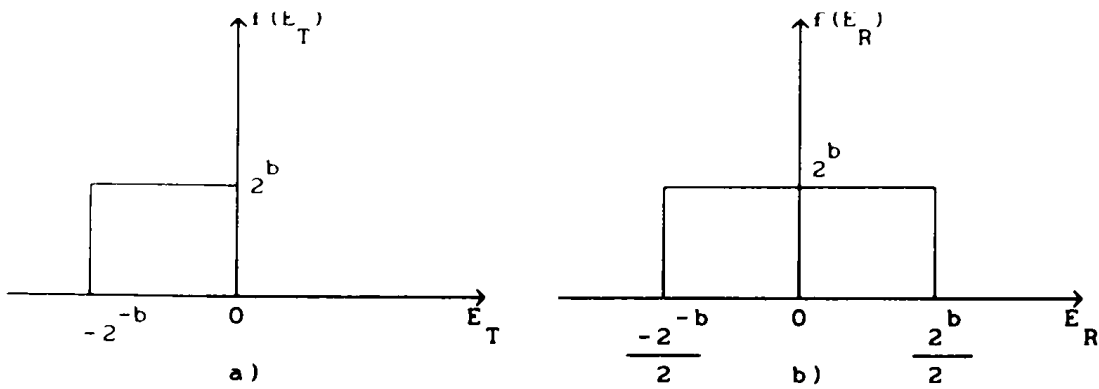


Figure 4.4 : Fonction densité de probabilité de l'erreur de quantification de la représentation en complément à 2.
 a) Par troncature b) Par arrondi

b) Arrondi

On procède de la même façon que pour la représentation en signe et valeur absolue et on trouve [30]:

$$\frac{-2^{-b}}{2} \leq E_R < \frac{2^{-b}}{2} \tag{4.22}$$

Pour toutes les valeurs de x. La variance du bruit d'arrondi est [30][31][7]:

$$\sigma_{E_R}^2 = \frac{2^{-b}}{12} \tag{4.23}$$

La caractéristique du quantificateur et la densité de probabilité de l'erreur E_R sont montrés dans les figures (4.3.b) et (4.4.b) respectivement.

4.2.3 Représentation en complément à 1 :

Dans cette représentation, on exprime les nombres comme suit [31]:

$$\begin{aligned} Q[x] &= (0_{\Delta} m_1 m_2 \dots m_b)_2, \quad x \geq 0 \\ &= (1_{\Delta} n_1 n_2 \dots n_b)_2, \quad x \leq 0 \end{aligned} \tag{4.24}$$

où

$$\begin{aligned} \text{pour } x \geq 0 \quad & (0_{\Delta} m_1 m_2 \dots m_b)_2 = |Q[x]| \\ \text{pour } x < 0 \quad & (0_{\Delta} n_1 n_2 \dots n_b)_2 = \text{Complément à 1 de } |Q[x]| \\ & = 1 - |Q[x]| - 2^{-b} \end{aligned}$$

Pour l'analyse de l'erreur introduite, on trouve les mêmes résultats que pour la représentation en signe et valeur absolue.

4.3 Effets de la longueur des mots sur les filtres IIR :

Les sources de l'erreur de quantification dans l'implantation des filtres numériques sont [30][33]:

- 1- La quantification des coefficients.
- 2- La quantification du signal d'entrée.
- 3- La quantification du produit des multiplications.
- 4- Les oscillations limites.
- 5- Les oscillations de dépassement.

4.3.1. Quantification du signal d'entrée :

Dans les systèmes de traitement numérique du signal, les signaux à traiter sont continus dans le temps et dans l'amplitude. Ils doivent, de ce fait, être discrétisés par une opération d'échantillonnage. Ensuite, vu la longueur finie des mots, ils doivent être quantifiés.

Le dispositif qui réalise la quantification est un convertisseur Analogique Numérique (CAN). On suppose que les données à la sortie du convertisseur sont représentées par une fraction à virgule fixe en complément à deux sur une longueur de $(b+1)$ bits. Les données à l'entrée sont arrondies au niveau de quantification le plus près pour obtenir la donnée quantifiée. L'écart entre deux valeurs quantifiées voisines est appelé pas de quantification et noté q . Il est donné par :

$$q = 2^{-b}$$

La figure (4.5) montre la caractéristique de quantification d'un convertisseur à $b = 2$

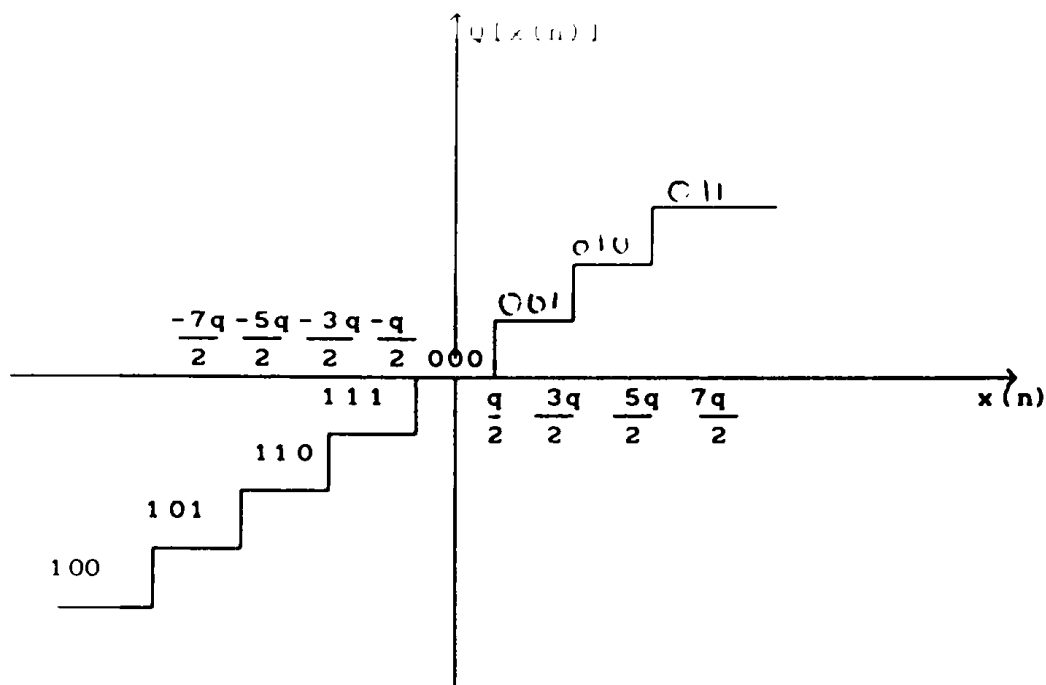


Figure 4.5 : Caractéristique d'un CAN pour $b = 2$

L'erreur de quantification est [7] :

$$e(n) = Q[x(n)] - x(n) \quad (4.25)$$

Cette erreur est bornée par l'intervalle [7] :

$$-\frac{q}{2} \leq e(n) \leq \frac{q}{2} \quad (4.26)$$

Pour analyser cette erreur, on utilise un modèle statistique qui sera adopté pour décrire l'effet de la quantification dans les algorithmes du traitement du signal [7] :

On admet ce qui suit [6]

1- La séquence d'erreur $\{e(n)\}$ est une séquence échantillon d'un processus aléatoire stationnaire.

2- La séquence d'erreur est non-corrélée avec la séquence des échantillons $\{x(n)\}$.

3- Les variables aléatoires du processus d'erreur sont non-corrélées. L'erreur est un bruit blanc.

4- La densité de probabilité du processus d'erreur est uniforme sur l'intervalle de l'erreur de quantification.

Sous ces considérations, le bruit de quantification est [7]:

$$\sigma_e^2 = \frac{q^2}{12} = \frac{2^{-2b}}{12} \quad (4.27)$$

Quand le signal quantifié est l'entrée d'un système linéaire invariant, la réponse peut être exprimée par :

$$y'(n) = y(n) + f(n) \quad (4.28)$$

où $y(n)$ est la réponse à $x(n)$ et $f(n)$ est la réponse à $e(n)$.

La moyenne et la variance du bruit à la sortie sont calculées en utilisant les équations (1.102) et (1.113) :

$$m_f = m_e \sum_{n=-\infty}^{+\infty} h(n) = m_e H(0) \quad (4.29)$$

$$\sigma_f^2 = \sigma_e^2 \sum_{n=-\infty}^{+\infty} |h(n)|^2 = \frac{\sigma_e^2}{2\pi} \int_{-\infty}^{+\infty} |H(\omega)|^2 d\omega$$

4.3.2 Effets de la quantification du produit des multiplications

Pour l'arithmétique à virgule fixe, le produit de deux nombres de longueur b -bits génère des nombres de longueur $2b$ bits. Pour pouvoir stocker ce résultat, on doit l'arrondir à b -bits. Ceci a pour effet qu'un bruit additif se superpose au signal à la sortie du mutiplieur. Le multiplieur peut être modélisé par la figure (4.6) et peut être exprimé par [6]:

$$Q[a_1 x(n)] = a_1 x(n) + e(n) \quad (4.30)$$

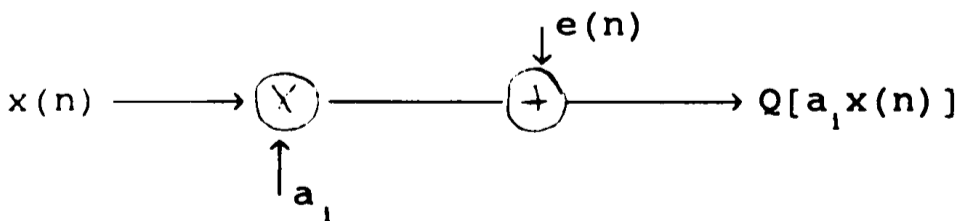


Figure 4.6 : Modele statistique du multiplieur

En utilisant le modèle du multiplieur, on peut schématiser la cellule du second ordre forme 1D selon la figure (4.7).

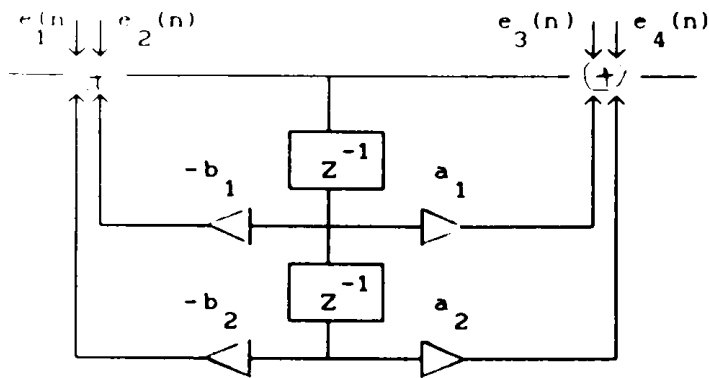


Figure 4.7 : Cellule du second ordre avec bruit de quantification des multiplieurs

Comme on le voit sur cette figure, pour chaque multiplieur, un signal d'erreur $e_i(n)$ est additionné au noeud de somme. Le signal d'erreur est modélisé par un processus aléatoire avec une densité de probabilité uniforme. La moyenne, la variance et la fonction d'autocorrélation sont données par :

$$E[e_i(n)] = 0 \quad (4.31)$$

$$\sigma_e^2 = \frac{q^2}{12} = \frac{2^{-2b}}{12} \quad (4.32)$$

$$\phi_e(m) = \frac{q^2}{12} \cdot \delta(m) \quad (4.33)$$

La densité spectrale de puissance est donnée par :

$$P_e(\omega) = \frac{q^2}{12} \quad (4.34)$$

Pour analyser le bruit d'arrondi des résultats des multiplieurs, il est nécessaire de définir la fonction de transfert propre de la source du bruit du filtre car les sources de bruit prennent place dans plusieurs endroits [32] [6].

Si $h_k(m)$ est la réponse impulsionnelle du filtre de la source du bruit à la sortie du filtre, alors la réponse à $e_k(n)$ est donnée par la convolution discrète [6]:

$$e_{ok}(n) = \sum_{m=0}^{\infty} h_k(m) e_k(n-m) \quad (4.35)$$

La variance de $e_{ok}(n)$ est donnée par [6]:

$$\sigma_{e_0}^2 = \sigma_e^2 \sum_{m=0}^{\infty} h_k^2(m) \quad (4.36)$$

La variance du bruit total est donnée par [6]:

$$\sigma_{e_0}^2 = \sum_{l=1}^M \sigma_{0l}^2 \quad (4.37)$$

L'évaluation de $\sigma_{e_0}^2$ nécessite le calcul de la somme de h_k^2 , cette somme est donnée par [6]:

$$\begin{aligned} \sum_{m=0}^{\infty} h_k^2(m) &= \frac{1}{2\pi j} \oint_C H_{k,M}(z) \cdot H_{k,M}(z^{-1}) \cdot z^{-1} dz \\ &= \frac{1}{2\pi} \int_0^{2\pi} |H_{k,M}(e^{j\omega})|^2 d\omega \end{aligned} \quad (4.38)$$

où $H_{k,M}(z)$ est la fonction de transfert du bruit de sa source à la sortie du filtre.

Il est clair que l'analyse du bruit dépend de la structure du filtre puisque la fonction de transfert du bruit varie d'une structure à une autre. Pour la structure 1D, on montre que la variance du bruit de sortie est donnée par [6]:

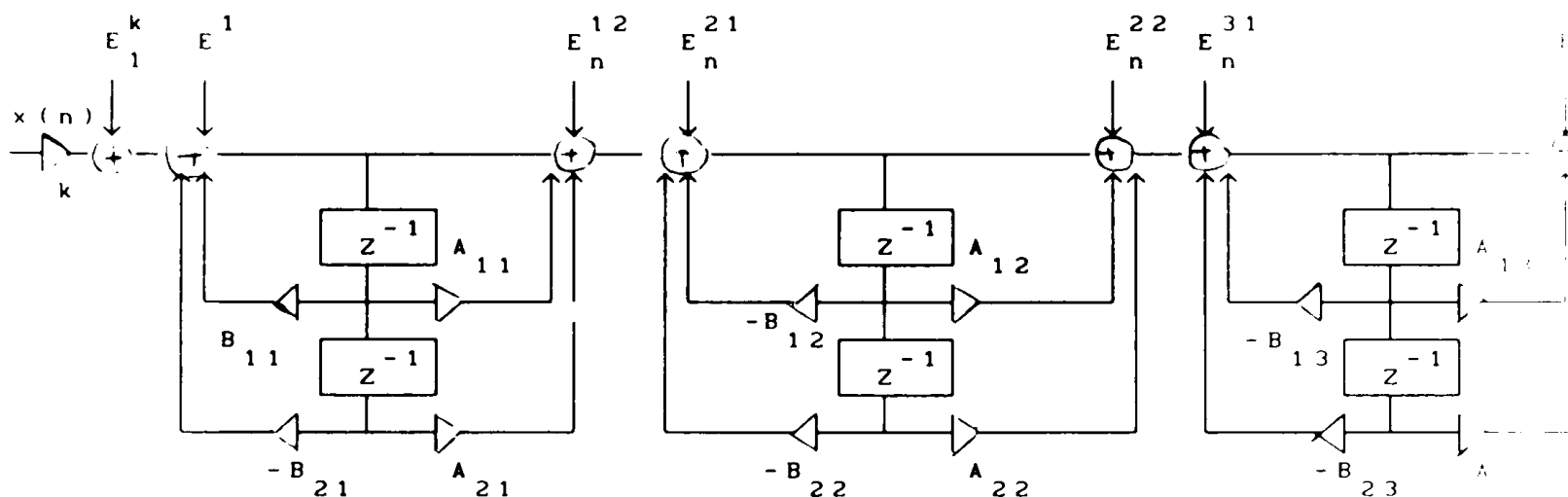
$$\begin{aligned} \sigma_{e_0}^2 &= \frac{\sigma_e^2}{2\pi} \left\{ r \sum_{k=1}^M \int_0^{2\pi} |H_{k,M}(e^{j\omega})|^2 d\omega + S \sum_{k=2}^M \int_0^{2\pi} |H_{k,M}(e^{j\omega})|^2 d\omega \right. \\ &\quad \left. + S + \int_0^{2\pi} |H_{1,M}(e^{j\omega})|^2 d\omega \right\} \end{aligned} \quad (4.39)$$

où σ_e^2 = variance du bruit d'arrondi

r = nombre de sources de bruits au premier noeud de sommation

- S = nombre de sources de bruit au second noeud de sommation
- M = nombre de cellules du second ordre
- K = index des sections de second ordre.

La figure (4.8) illustre un exemple de la réponse au bruit de trois sections de second ordre mises en cascade.



$$\begin{aligned}
 Y(z) &= X(z) \cdot K \cdot H_1(z) \cdot H_2(z) \cdot H_3(z) = \text{sortie du filtre idéal} \\
 &+ E_1^k(z) \cdot H_{1,3}(z) = \text{bruit de sortie du coefficient K} \\
 &+ E_1^{11}(z) \cdot H_{1,3}(z) + E_n^{12}(z) \cdot H_{2,3}(z) = \text{bruit de sortie de la section 1} \\
 &+ E_n^{21}(z) \cdot H_{2,3}(z) + E_n^{22}(z) \cdot H_{3,3}(z) = \text{bruit de sortie de la section 2} \\
 &+ E_n^{31}(z) \cdot H_{3,3}(z) + E_n^{32}(z) = \text{bruit de sortie de la section 3}
 \end{aligned}$$

Figure 4.8 : Reponse du bruit d'arrondi de trois cellules de second ordre mises en cascade.

La relation (4.39) peut être réécrite pour d'autres structures.

La fonction de transfert du bruit est donnée par :

$$H_{k,M}(z) = H_k(z) \cdot H_{k+1}(z) \dots H_M(z) \quad (4.40)$$

La fonction de transfert du bruit dépend de l'organisation des cellules. Pour M cellules, il y a $(M!)^2$ organisations possibles. Pour chacune d'elles, on évalue l'équation (4.39) et on garde celle qui minimise cette équation. Ceci exige un grand nombre d'évaluations. Par exemple, pour M=4, on a 576 évaluations et pour M=5, on 14400 évaluations. On montre [6] que l'appariement optimal des pôles et des zéros peut être obtenu en couplant chaque pôle avec le zéro qui lui est le plus près dans le plan Z. La figure (4.11) en montre un exemple. Ceci réduit le nombre d'évaluations

nécessaires de l'équation (4.39) à (M!).

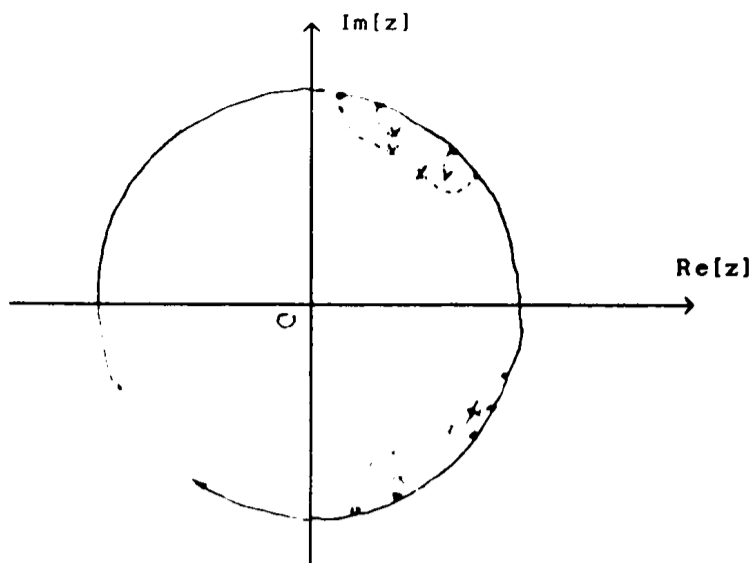


Figure 4.9 : Couplement des poles et des zeros

Dynamique du filtre :

La dynamique du filtre est définie comme étant le nombre de bits non affectés par le bruit. La dynamique peut être exprimée en décibel par [6]:

$$DR = (\text{BIT}_N - 1) 20 \text{ Log}_{10} 2 \quad (4.41)$$

où BIT_N est la position du bit le moins significatif du bruit calculé. Il est donné par [6]:

$$\text{BIT}_N = \text{INT} \left(\frac{-\text{Log } \sigma_{e0}}{\text{Log } 2} \right) \quad (4.42)$$

où INT dénote la partie entière.

4.3.4 Les dépassements et mise à l'échelle :

Dans la réalisation des filtres numériques, il convient de transmettre un signal dont la dynamique est aussi élevée que possible : un signal dont l'amplitude maximale soit la plus grande que possible par rapport au pas de quantification. Cette condition doit être satisfaite en présence d'une contrainte : en aucun point du filtre, le signal ne peut exéder la capacité des registres utilisés pour le contenir [11].

Les points où cette contrainte de non dépassement doit être

vérifiées sont les entrées des multiplieurs. En effet, les résultats partiels de plusieurs sommations successives peuvent excéder la capacité de l'accumulateur même si la somme totale est inférieure à cette capacité [11].

En virgule fixe, tout dépassement de capacité est interprété comme un changement de signe : lorsque ce phénomène se produit à la sortie d'un additionneur, il peut entraîner une oscillation forcée permanente de grande amplitude. Cette oscillation subsiste même en absence du signal d'entrée [33]. On montre que ces oscillations peuvent être annulées en utilisant des additionneurs à logique saturée.

La figure (4.10) montre la structure cascade avec les points où la contrainte de non dépassement doit être vérifiée, notée par (*).

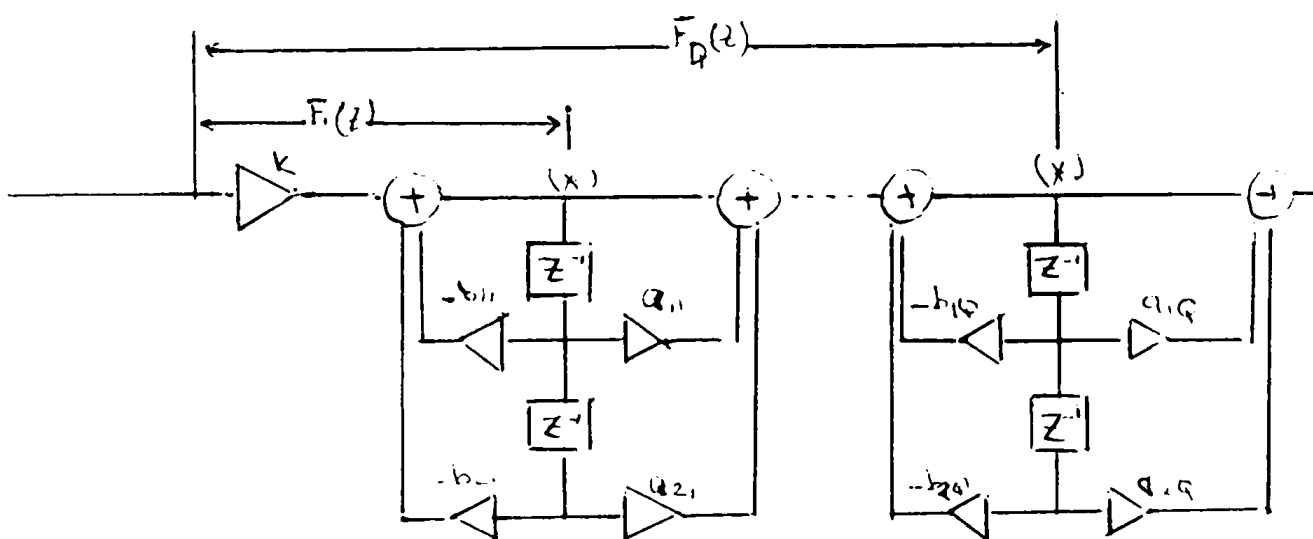


Figure 4.10 : Points où le non dépassement doit être vérifié pour la structure cascade.

Une mise à l'échelle du filtre numérique doit être opérée pour que ces débordements n'aient pas lieu.

Soient $F_1(z)$ la fonction de transfert de l'entrée au point de dépassement avant la mise à l'échelle et $F'_1(z)$ la fonction de

transfert après la mise à l'échelle.

Si $H(z)$:

$$H(z) = k \prod_{i=1}^0 \frac{1 + a_{1i} z^{-1} + a_{2i} z^{-2}}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \quad (4.43)$$

dénote la fonction de transfert du filtre, alors la version mise à l'échelle aura pour fonction de transfert [34]:

$$H(z) = k' \prod_{i=1}^0 \frac{\alpha_{0i} + \alpha_{1i} z^{-1} + \alpha_{2i} z^{-2}}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \quad (4.44)$$

Les fonctions $F_i(z)$ et $F'_i(z)$ pour la forme 1D sont données par [34]:

$$F_i(z) = \frac{K}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \prod_{j=1}^{i-1} \frac{1 + a_{1j} z^{-1} + a_{2j} z^{-2}}{1 + b_{1j} z^{-1} + b_{2j} z^{-2}} \quad (4.45)$$

$$F'_i(z) = \frac{K'}{1 + b_{1i} z^{-1} + b_{2i} z^{-2}} \prod_{j=1}^0 \frac{\alpha_{0j} + \alpha_{1j} z^{-1} + \alpha_{2j} z^{-2}}{1 + b_{1j} z^{-1} + b_{2j} z^{-2}} \quad (4.46)$$

avec $\prod_{j=1}^0 (\cdot) = 1$

On montre [11] que si $X(\omega)$, la transformée de Fourier du signal d'entrée, est telle que $\|X(\omega)\|_q \leq M$, pour $q \geq 1$, on en déduit que $\|X_n\|_q \leq M$ et pour que le débordement n'ait pas lieu, il suffit d'avoir :

$$\|F'_i\| \leq 1, \quad \|X\|_q \leq M; \quad \frac{1}{p} + \frac{1}{q} = 1 \quad (4.47)$$

p et q sont des entiers positifs et $\|\cdot\|_p$ dénote la norme L_p .

La mise à l'échelle doit être faite comme suit [34]:

$$K' = S_1$$

$$\alpha_{1j} = \frac{S_{1+1}}{S_1} a_{1j} ; i = 0, 1, 2 , j = 1 \dots Q \quad (4.48)$$

$$S_{0+1} = K$$

$$S_1 = 1/||F_1||_p$$

De ceci, on a :

$$F'_1(z) = S_1 \cdot F_1(z) \quad (4.49)$$

Similairement, on tire pour la forme 2D [34]:

$$F_1(z) = k \prod_{j=1}^{i=1} \frac{1 + a_{11} z^{-1} + a_{21} z^{-2}}{1 + b_{11} z^{-1} + b_{21} z^{-2}} \quad (4.50)$$

$$\text{et} \quad F'_1(z) = \prod_{j=1}^{i-1} \frac{\alpha_{0j} + \alpha_{1j} z^{-1} + \alpha_{2j} z^{-2}}{1 + b_{1j} z^{-1} + b_{2j} z^{-2}} \quad (4.51)$$

La fonction de transfert H(z) est de la forme (4.44).

La mise à l'échelle est faite comme suit [34]:

$$k' = \frac{k}{S_H}$$

$$\alpha_{1j} = \frac{S_1}{S_{1-1}} a_{1j} ; i=0, 1, 2 ; j=1 \dots Q \quad (4.52)$$

$$S_0 = 1$$

$$\text{et} \quad S_1 = 1/||F_1||_p$$

La figure (4.11) montre la structure cascade mise à l'échelle pour les formes 1D et 2D.

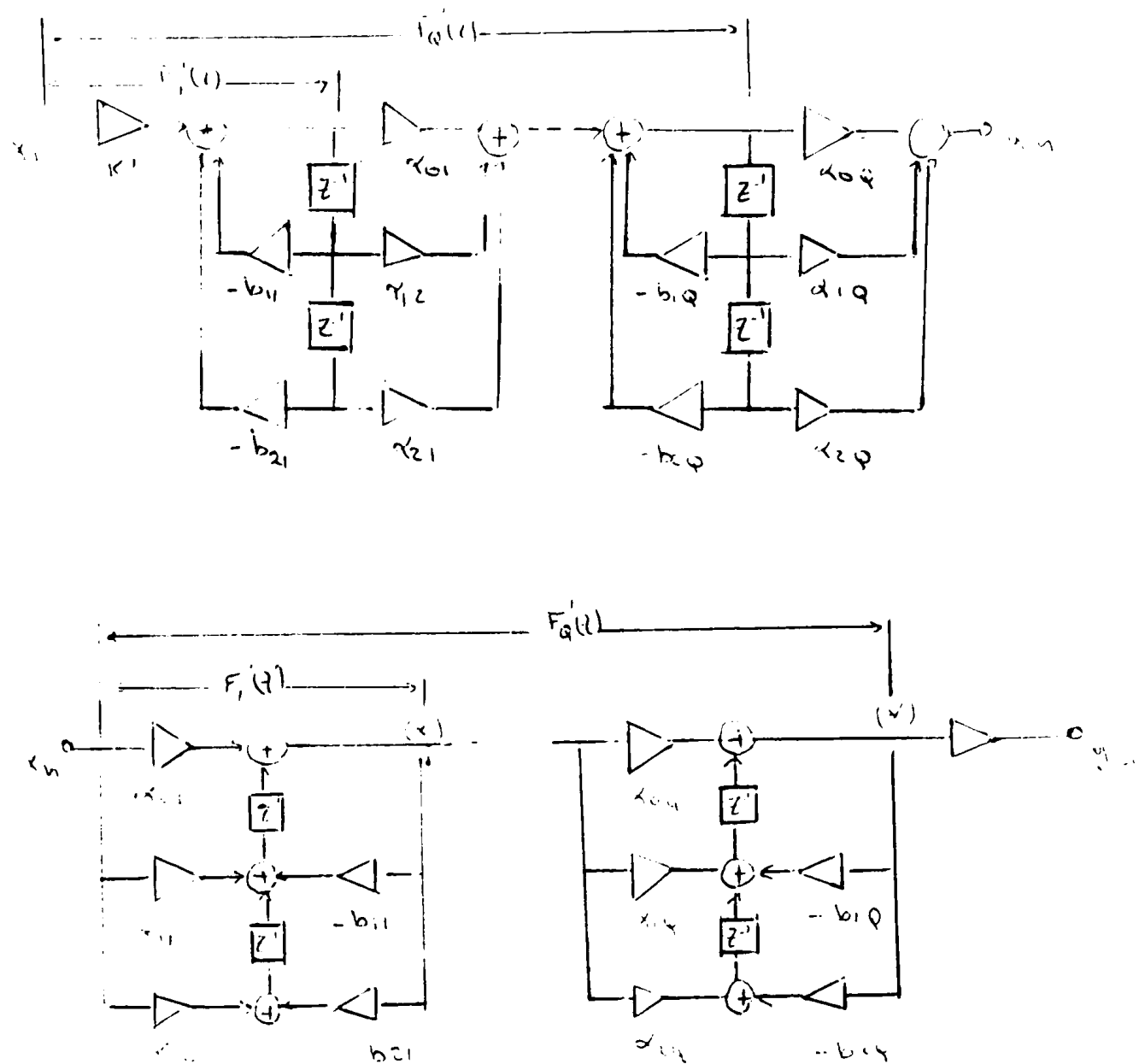


Figure 4.11 : Structure cascade mise a l'echelle
 a) Forme 1D b) forme 2D

Du point de vue pratique, les valeurs les plus significatives de p et q dans la relation (4.47) sont 1, 2 et ∞ :

- Le cas $p = 1, q = \infty$ impose que $X(\omega)$ soit partout bornée par M
- Le cas $p = q = 2$ est appliquée pour les signaux à énergie finie
 en effet on a :

$$\|X\|_2^2 = \text{Energie du signal.}$$

- Le cas $p = \infty$ et $q = 1$ conduit à la condition la plus sévère sur F_1' car on a, toujours :

$$\|F_1'\|_p \leq \|F_1'\|_\infty \quad \forall p \geq 1$$

4.3.5 Les cycles limites :

Quand on excite l'entrée d'un filtre IIR par une séquence $x(n)$ qui est constante pour un certain instant :

$$x(n) = \begin{cases} A & \text{pour } n = n_0 \\ 0 & \text{ailleurs} \end{cases}$$

la sortie doit tendre idéalement vers zéro. Cependant quand on implante le filtre par des registres à longueur finie, la sortie peut osciller dans un intervalle d'amplitude non nulle. Ce phénomène est appelé "oscillation à cycles limites" [6]. Il importe de connaître une borne supérieure de ces oscillations.

Une analyse déterministe de ce phénomène est très difficile et ne peut être obtenue que par simulation [11].

La table 4.1 illustre ce phénomène pour le filtre du premier ordre :

$$y(n) = x(n) - 0.5 y(n-1)$$

pour une entrée :

$$x(n) = \begin{cases} 0.875 & \text{pour } n = 0 \\ 0 & \text{ailleurs} \end{cases}$$

en supposant que les données sont représentées sur 3 bits :

n	idéal	quantifié	Représentation binaire
0	0.875	0.875	0.111
1	-0.4375	-0.5	1.100
2	0.21875	0.25	0.010
3	0.0546875	-0.125	1.111
4	0.02734375	0.125	0.001
∞	0.0000	±0.125	1.111 ou 0.001

Table 4.1 : Exemple de cycles limites

On remarque qu'au lieu d'être zéro, la sortie oscille entre 0.125 et -0.125

4.4 Effets de la longueur des mots sur la FFT :

Le graphe de fluence de l'algorithme de la FFT à entrelacement temporel est montré sur la figure (4.12) pour N=8.

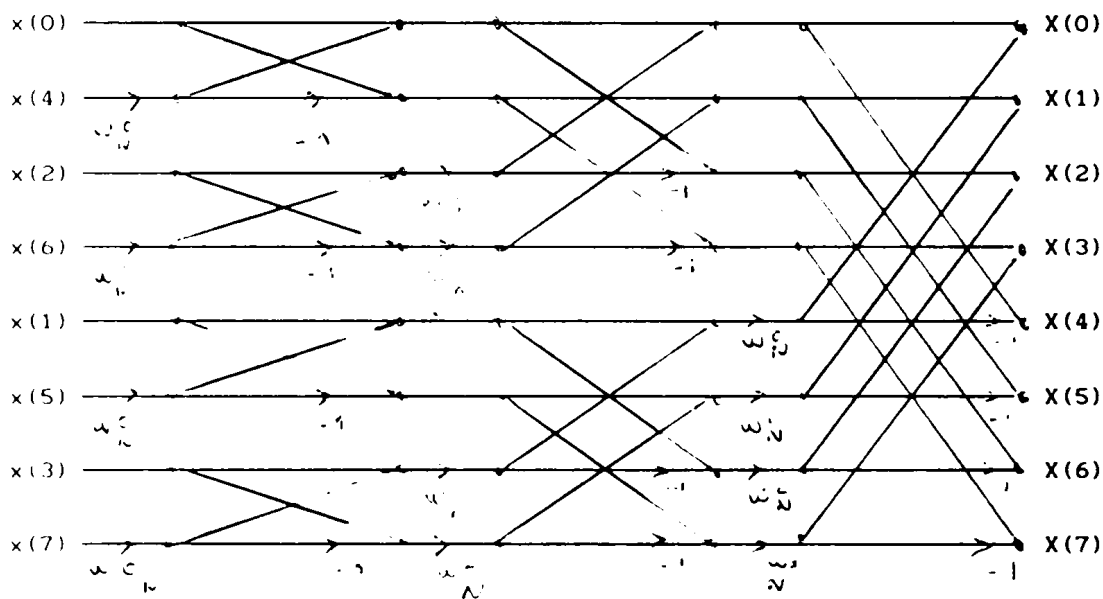


Figure 4.12 : Graphe de fluence de la FFT à entrelacement temporel

L'opération de base de l'algorithme est le papillon décrit par les équations :

$$\begin{aligned}
 X_{m+1}(p) &= X_m(p) + W_N^r X_m(q) \\
 X_{m+1}(q) &= X_m(p) - W_N^r X_m(q)
 \end{aligned}
 \tag{4.53}$$

L'indice m dénote l'étage m , l'indice $(m+1)$ dénote l'étage $(m+1)$ et les indices p et q dénotent la position des échantillons dans l'étage. Le papillon est montré sur la figure 4.15.a

Les produits $X_m(q) \cdot W_N^r$ doivent être quantifiés pour une implantation par une arithmétique à précision finie. Ceci introduit une erreur $e(m,q)$ qui se superpose au signal de sortie. Cette erreur est une séquence complexe.

La figure (4.13.b) montre le modèle du papillon avec l'erreur introduite.

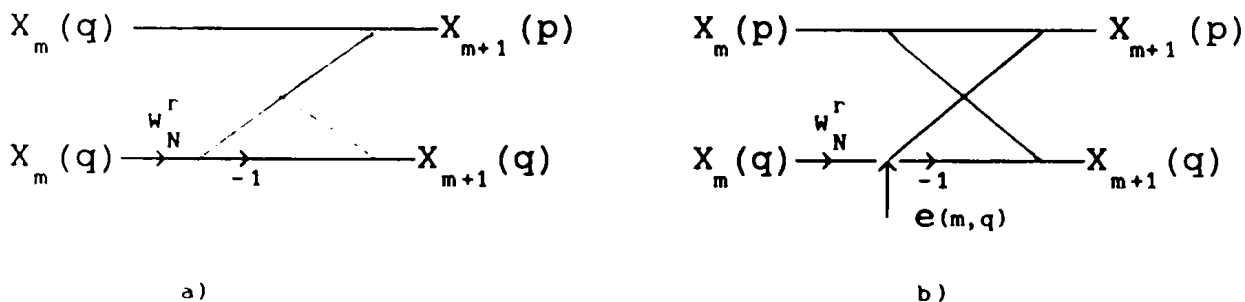


Figure 4.13 : a) Graphe de fluence d'un papillon
b) Graphe de fluence d'un papillon avec l'erreur d'arrondi

Si on suppose que chaque multiplication complexe est réalisée par quatre multiplications réelles et que chaque multiplication réelle est un bruit blanc dont la variance est $(2^{-2b}/12)$ [6,7,31], la variance de $e(m,q)$ sera :

$$\sigma_B^2 = \frac{2^{-2b}}{12} \quad (4.54)$$

b étant le nombre de bits utilisés pour représenter la partie fractionnaire.

On remarque à partir de la figure (4.12) que :

1- Le bruit généré par un papillon à la sortie d'un étage donné se propage à la sortie de l'étage suivant multiplié par une constante d'amplitude unité (-1 ou W_N^r). De ce fait, sa variance restera inchangée en atteignant la sortie de l'étage final.

2- Chaque point de sortie $X(k)$ est relié à :

- Un papillon de l'étage final M-1
- 2 papillons de l'étage M-2
- 4 papillons de l'étage M-3
-
-
-
- 2^m papillons de l'étage M -(m-1)

avec $N = 2^M$ et $m = 0 \dots M-1$

Donc chaque échantillon de sortie $X(k)$ de la TFD se trouve lié à $(N-1)$ papillons. L'erreur $E(k)$ qui lui est entachée a une variance:

$$\sigma_E^2 = E[|E(k)|^2] = (N-1)\sigma_B^2 \quad (4.55)$$

qui donnera pour N assez grand

$$\sigma_E^2 = N.\sigma_B^2 \quad (4.56)$$

Pour qu'il n'y ait pas de débordement à la fin des calculs, une analyse doit être faite. On a à partir de l'équation (4.53):

$$\max [|X_{m+1}|] \leq 2.\max[|X_m|] \quad (4.57)$$

Cette relation implique que le module maximum peut doubler d'un étage à un autre. Ce qui donne :

$$\max[|X(k)|] \leq N.\max[|x(n)|] \quad (4.58)$$

Dans la représentation en virgule fixe, un nombre X est tel que :

$$-1 \leq X \leq 1$$

ou $|X| \leq 1$

Donc, pour qu'il n'y ait pas débordement, on doit avoir :

$$|x(n)| < \frac{1}{N} \quad , \quad 0 \leq n \leq N-1 \quad (4.59)$$

pour garantir

$$|X(k)| < 1 \quad , \quad 0 \leq k \leq N-1 \quad (4.60)$$

Si on suppose que $x(n)$ est un signal blanc dont la partie réelle et imaginaire sont uniformément réparties entre $(-1/\sqrt{2}N)$ et $(1/\sqrt{2}N)$, alors sa variance sera :

$$\sigma_x^2 = \frac{1}{3N^2} \quad (4.61)$$

et celle de $X(k)$ est :

$$\sigma_x^2 = \frac{1}{3N} = N \cdot \sigma_x^2 \quad (4.62)$$

Le rapport bruit sur signal est défini par :

$$\text{NSR} = \frac{\sigma_E^2}{\sigma_x^2} = \frac{\text{Puissance du bruit}}{\text{Puissance du signal}} \quad (4.63)$$

Pour le cas considéré, on aura :

$$\text{NSR} = 3N^2 \sigma_B^2 = N^2 2^{-2b} \quad (4.64)$$

L'interprétation de ce résultat est que le rapport bruit sur signal augmente de un bit par étage [7,31].

Ce résultat est médiocre pour les applications où la précision de calcul est exigée. Cette méthode de prévention contre le débordement est appelée Mise à l'échelle des données d'entrée.

Une alternative à cette méthode consiste à insérer des atténuations de 1/2 à l'entrée de chaque étage. A l'entrée du premier étage, les données sont telles que :

$$|x(n)| < 1$$

La figure (4.14) montre le modèle du papillon avec une atténuation de 1/2 à l'entrée de chaque étage.

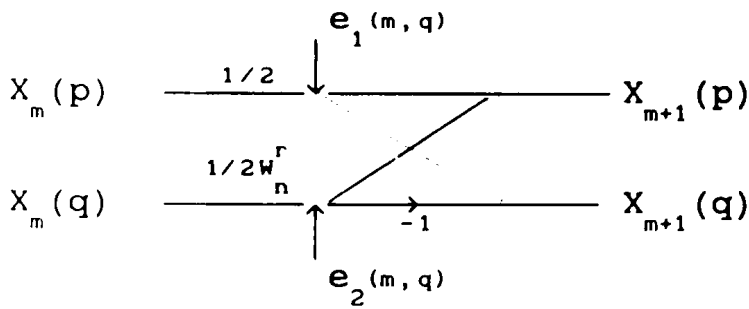


Figure 4.13 : Graphe de fluence d'un papillon avec bruit d'arrondi et atténuation de 1/2

Dans ces deux cas, la sortie du dernier étage ne consiste pas en la TFD du signal mais en 1/N fois cette TFD. Mais la différence est que la seconde méthode introduit beaucoup moins de bruit.

En effet, l'analyse du modèle de la figure (4.14) révèle qu'on a quatre multiplications qui donnent un bruit de variance σ_b^2 (équation 4.54) et si les atténuations sont faites par arrondi, elles introduisent un bruit de variance [6]:

$$\sigma_s^2 = \frac{2^{-2b}}{2} \quad (4.65)$$

Ce qui donne un bruit total [6]

$$\sigma_{BS}^2 = \frac{5}{6} 2^{-2b} \quad (4.66)$$

Ce bruit, en se propageant d'un étage à un autre, se trouve divisé par 2. Ce qui donne qu'un bruit généré au $m^{1\text{eme}}$ étage se trouve à la sortie du dernier étage multiplié par une constante $(1/2)^{M-m-1}$.

En sachant que chaque point de sortie est lié à 2^{M-m-1} papillons du $m^{1\text{eme}}$ étage, soit 2^{M-m-1} sources de bruit généré au $m^{1\text{eme}}$ étage

($m = 0 \dots M-1$)

Ceci donne un bruit de sortie $E(k)$ de variance :

$$\begin{aligned}
\sigma_E^2 &= \sigma_{BS}^2 \sum_{m=0}^{M-1} 2^{M-m-1} \left(\frac{1}{2}\right)^{2m} \\
&= \sigma_{BS}^2 \sum_{m=0}^{M-1} \left(\frac{1}{2}\right)^{M-m-2} \\
&= \sigma_{BS}^2 \sum_{l=0}^{M-1} \left(\frac{1}{2}\right)^l \\
&= 2 \sigma_{BS}^2 \left(1 - \left(\frac{1}{2}\right)^M\right)
\end{aligned}
\tag{4.67}$$

Pour un nombre d'échantillons assez grand, on obtient :

$$\sigma_E^2 = 2 \sigma_{BS}^2 = \frac{5}{3} 2^{-2b}
\tag{4.68}$$

En supposant qu'à l'entrée, on a un signal blanc comme précédemment, on trouve un rapport bruit :

$$NSR = \frac{\sigma_E^2}{\sigma_X^2} = 3.N \sigma_{BS}^2 = 5.N2^{-2b}
\tag{4.69}$$

Le rapport bruit sur signal augmente de un bit par étage ce qui constitue une grande amélioration du résultat précédent. Cette méthode s'appelle "Mise à l'échelle inconditionnelle des données"

Une troisième méthode de prévention contre le débordement consiste à utiliser le format flottant par bloc. Dans cette procédure, les données à l'entrée sont normalisées au plus loin vers la gauche du registre avec la restriction que $|x(n)| < 1$. Le calcul commence et à chaque fin d'étage, on teste si un débordement a lieu, si c'est le cas, on décale toutes les données de un bit vers la droite et le calcul continue. Le nombre de décalages nécessaires est calculé pour déterminer l'exposant commun à toutes les données.

La variance du bruit total ainsi que le rapport bruit sur signal dépendent du nombre de décalages fait ainsi que de l'étage où ils ont eu lieu. Cette méthode s'appelle "Mise à l'échelle

conditionnelle des données". Le rapport bruit sur signal sera inférieur à celui des deux premières méthodes et à la limite, il sera égal à celui du second quand les décalages sont faits dans tous les étages.

4.5 Effets de la longueur des mots sur les filtres FIR :

La structure la plus utilisée pour réaliser des filtres FIR est la structure directe (voir chapitre 3). La longueur finie des mots a un effet sur les coefficients des filtres ce qui pose le problème de la sensibilité (voir chapitre 3), sur le résultat des multiplications et sur le débordement.

4.5-2 Effets de la quantification des résultats des multiplications

Pour le filtre FIR, deux cas sont à considérer :

1^{er} cas : On arrondit le résultat de chaque multiplication et puisqu'il y a N multiplications, le bruit d'arrondi total a pour variance [7]:

$$\sigma_B^2 = N \cdot \frac{q^2}{12} = N \frac{2^{-2b}}{12} \quad (4.70)$$

2^{ème} cas : On n'arrondit le résultat qu'après avoir fait l'accumulation de tous les produits. Dans ce cas, une seule source d'arrondi est présente. Le bruit d'arrondi a une variance de [6]:

$$\sigma_B^2 = \frac{q^2}{12} = \frac{2^{-2b}}{12} \quad (4.71)$$

Ce deuxième cas est le plus fréquent dans l'implantation des applications. Tous les processeurs de signal disposent d'un multiplieur accumulateur qui multiplie deux nombres de b-bits et accumule les résultats sur 2b-bits.

4.5-3 Prévention contre le débordement

Pour prévenir contre le débordement, le signal à l'entrée doit être mis à l'échelle par une constante C telle que [6]:

$$C = \frac{1}{\sum_{n=0}^{N-1} |h(n)|} \quad (4.72)$$

SUR UN PROCESSEUR DE SIGNAL

5.1 Introduction :

Les algorithmes de base du traitement numérique du signal sont la transformée de Fourier rapide (FFT), le filtrage numérique (IIR et FIR) et le produit de convolution. [1,3 et 4]. En se basant sur l'étude théorique des chapitres précédents, on a implémenté ces algorithmes sur le système de développement du processeur de la compagnie ANALOG DEVICES l'ADSP-2100. Le processeur est optimisé pour le traitement du signal et les applications en temps réel nécessitant des vitesses de calcul importantes. Il traite des mots de 16 bits en virgule fixe et fournit des possibilités de traitement en multiprécision. (Des mots de plus de 16 bits). Les algorithmes ayant été considérés traitent des données représentées en simple précision.

L'annexe 1 décrit le processeur et son système de développement.

Les applications en temps réel doivent s'exécuter dans le minimum de temps, générer le moins de bruit et offrir la plus grande dynamique du signal [5]. Les deux derniers points sont généralement exprimés par le rapport signal sur bruit (S/B). En pratique, il y a une dualité entre le temps de calcul et le rapport (S/B) et on est appelé à faire un compromis entre les deux selon les exigences.

Dans ce qui suit, on traitera les algorithmes implantés en indiquant leurs performances et éventuellement leurs limites. Le produit de convolution sera considéré avec le filtrage FIR. Le compromis cité ci-dessus sera mis en évidence dans le cas de la FFT.

Des signaux modèles permettent de vérifier le bon fonctionnement des algorithmes. Un de ces signaux montrera la limite de la représentation en virgule fixe. Deux applications seront considérées : le filtrage d'un signal EEG contenant un

bruit additif et l'estimation de la fonction d'autocorrélation d'un signal aléatoire.

Les listings des programmes en langage assembleur de l'ADSP-2100 se trouvent dans l'annexe 2.

5.2 La FFT :

D'après le paragraphe 2.7.1, une FFT sur N points se calcule en M étages avec $M = \text{Log}_2 N$, chaque étage contient des groupes et chaque groupe contient des papillons.

La figure (5.1) montre l'organigramme de calcul de la FFT. Celle-ci se calcule en trois sous-routines :

- La première place les données dans l'ordre bit inversé.
- La deuxième calcule les échantillons de sortie d'un étage.
- La troisième met à l'échelle les données de sortie d'un étage.

Le calcul proprement dit se fait en trois boucles :

- Une boucle calcule les étages.
- Une boucle calcule les groupes.
- Une boucle calcule les papillons d'un groupe.

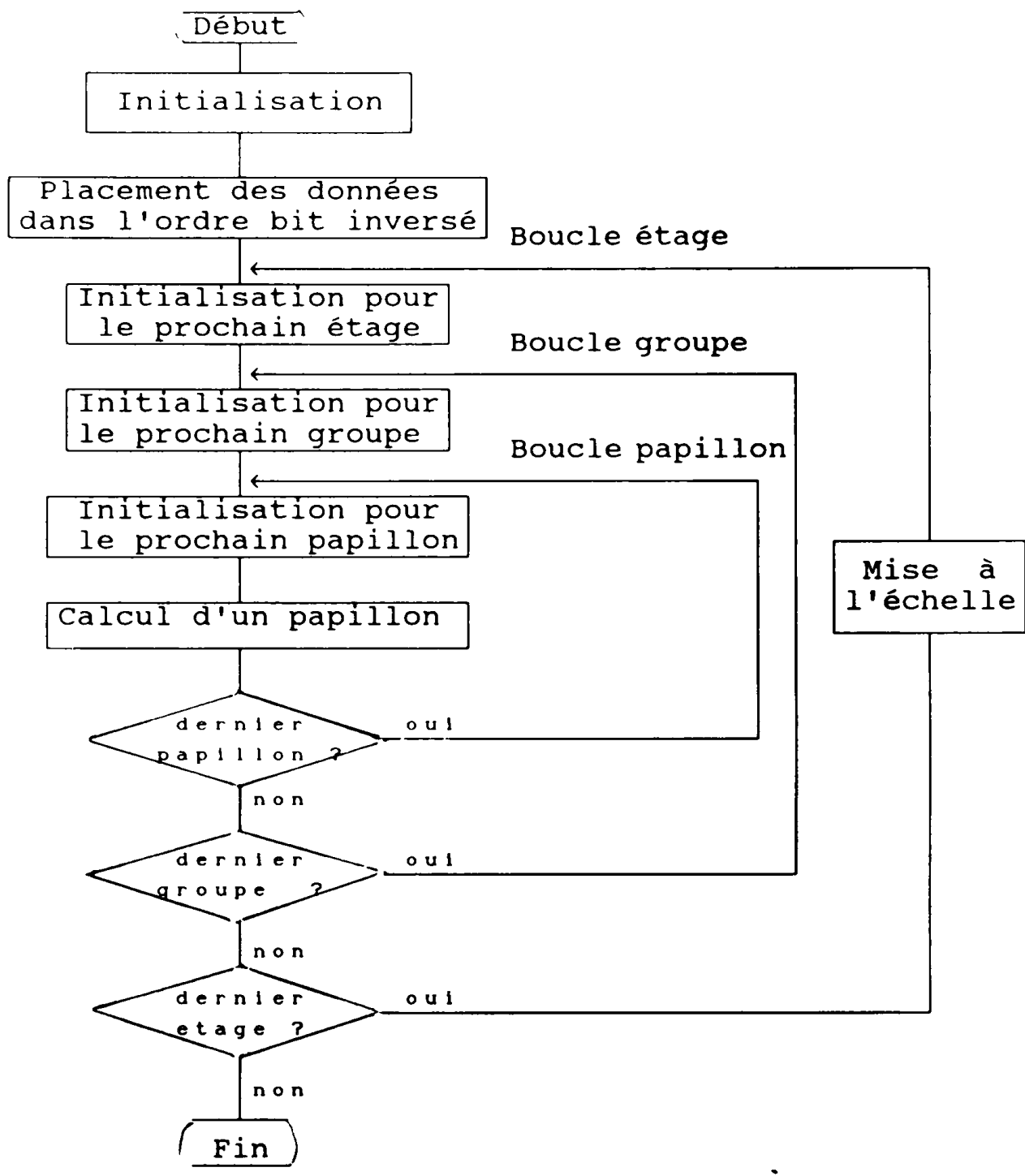


Figure 5.1 : Organigramme de calcul de la FFT

L'élément de calcul élémentaire de la FFT est le papillon. Son graphe de fluence est montré dans la figure (5.2). Les variables x et y représentent les parties réelles et imaginaires des données. L'exponentielle complexe est exprimée en partie réelle et imaginaire:

$$W_N = e^{-j2\pi/N} = \cos 2\pi/N - j\sin 2\pi/N = C + j(-S)$$

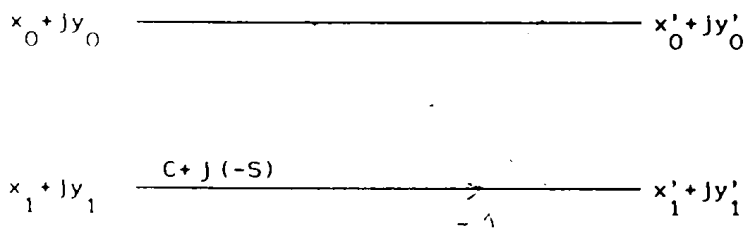


Figure 5.2 : Graphe de fluence d'un papillon

Le noeud dual est multiplié par le facteur $[C+j(-S)]$. Le résultat est additionné au noeud primaire pour produire $x'_0 + jy'_0$ et soustrait du noeud pour produire $x'_1 + jy'_1$. Les équations (5.1) à (5.4) calculent les parties réelles et imaginaires de la sortie des papillons.

$$x'_0 = x_0 + [C.x_1 - (-S)y_1] \quad (5.1)$$

$$y'_0 = y_0 + [C.y_1 + (-S)x_1] \quad (5.2)$$

$$x'_1 = x_0 - [C.x_1 - (-S)y_1] \quad (5.3)$$

$$y'_1 = y_0 - [C.y_1 + (-S)x_1] \quad (5.4)$$

Comme on l'a mentionné au quatrième chapitre , la dynamique du signal à la sortie de la FFT peut augmenter par un maximum de M bits. Ceci peut conduire à des débordements dans les calculs et à des résultats erronés. Il existe trois méthodes de prévention contre le débordement :

- La mise conditionnelle en format flottant par bloc (CBFP).
- La mise incondionnelle en format flottant par bloc (UBFP).
- La mise à l'échelle des données à l'entrée (INPUT).

Ces méthodes sont considérées ci-dessous.

5.2.1 La mise conditionnelle en format flottant par bloc

Dans cette méthode, un exposant commun est affecté aux données réelles et imaginaires constituant le signal d'entrée. A la fin du calcul de chaque étage, on teste si un débordement a eu lieu. Si tel est le cas, on met à l'échelle les données, on met à

jour l'exposant commun et on passe à l'étage suivant. Sinon, on passe directement à l'étage suivant. Le programme correspondant se trouve dans le paragraphe A2.1.1 de l'annexe 2.

Au cours du développement de ce programme, on a constaté les deux points suivants :

- Les deux premiers étages consomment un temps de calcul énorme. Pour une FFT sur 1024 points, cette consommation est de 25 % du temps de calcul total.

- Les papillons des deux premiers étages peuvent être calculés sans multiplications. Ils ont pour facteurs exponentiels (W_N^k) respectifs 1 et, 1 et -j. Les équations (5.1) à (5.4) deviennent :

pour $W_N^k = 1$

$$x'_0 = x_0 + x_1 \quad (5.5)$$

$$y'_0 = y_0 + y_1 \quad (5.6)$$

$$x'_1 = x_0 - x_1 \quad (5.7)$$

$$y'_1 = y_0 - y_1 \quad (5.8)$$

et pour $W_N^k = -j$

$$x'_0 = x_0 + y_1 \quad (5.9)$$

$$y'_0 = y_0 - x_1 \quad (5.10)$$

$$x'_1 = x_0 - y_1 \quad (5.11)$$

$$y'_1 = y_0 + x_1 \quad (5.12)$$

Le codage de ces équations en langage assembleur ferait réduire le calcul du papillon de 4 cycles processeur par rapport aux équations (5.1) à (5.4).

La prise en compte de ces deux points nous a conduit au calcul des deux premiers étages à part. On a gagné ainsi (12.83%) du temps de calcul pour $N = 1024$. Le gain sera plus significatif pour un nombre de points plus petit. Par exemple pour $N = 64$, le gain est de (19 %).

Du point de vue bruit de calcul, on note ce qui suit :

- Les deux premiers étages contiennent 75 % ($3N/4$) des sources de bruit présentes à la sortie de la FFT. La réalisation de ces étages sans multiplications éliminerait ces sources et réduirait celles qui sont présentes à la sortie à $(N/4 - 1)$ sources.
- L'ADSP-2100 contient un multiplieur/accumulateur dont l'exploitation fait réduire le nombre de sources de bruit par papillons à deux.

Pour éviter le calcul des facteurs exponentiels durant l'exécution du programme, ceux-ci sont initialement calculés et stockés dans une table dans la mémoire.

Dynamique du signal

En examinant les équations (5.1) à (5.4), on constate que les données d'entrée d'un papillon peuvent avoir leurs valeurs multipliées par $(1 + \sqrt{2})$, leur dynamique peut ainsi augmenter de 2 bits.

Pour cette raison, deux bits de garde sont initialement prévus pour que les débordements ne causent pas d'erreurs. Ceci réduit la dynamique du signal de 15 bits (nombre maximal de bits réservés par le processeur pour la partie fractionnaire) à 13 bits

Temps de calcul

Les papillons des deux premiers étages se calculent en huit cycles et ceux des autres étages se calculent en douze cycles.

Le temps de calcul d'une FFT sur N points dépend du nombre de débordements survenants au cours des calculs. Le temps minimal correspond à une FFT qui ne produit aucun débordement, ce temps est :

$$T_{\min} = (6N \cdot \text{Log}_2 N + 1.5N + 36 \text{Log}_2 N + 1) T_{\text{cycle}}$$

T_{cycle} étant le temps de cycle du processeur. Le temps maximal correspond à une FFT qui produit un débordement à chaque étage, ce temps est :

$$T_{\max} = (10N \cdot \text{Log}_2 N - 2.5N + 49M + 8) T_{\text{cycle}}$$

Si on a $T_{\text{cycle}} = 125 \text{ ns}$ et $N = 1024$, la FFT se calcule entre un

temps minimal, $T_{\min} = 7.92 \text{ ms}$ et un temps maximal, $T_{\max} = 12.54 \text{ ms}$

Bruit de calcul

Le calcul exact du bruit à la sortie dépend du nombre de débordements et de leur étage d'origine. Il n'y a pas de formule exacte de calcul de bruit.

APPLICATIONS

On a appliqué les programmes développés à quatre signaux modèles de longueur $N = 1024$. Ces signaux sont :

- Un signal Cosinus $\text{Sig1}(i) = \text{Cos}(2\pi i/N)$ $i=0\dots1023$

- Trois signaux constitués par la somme de 40 sinusoides d'amplitudes variées.

$$\text{Sig2}(i) = \sum_{k=0}^{40} 2^{-1/40} \text{Cos}(20\pi ik/N) \quad i = 0\dots1023$$

$$\text{Sig3}(i) = \sum_{k=0}^{40} 2^{-1/40} \text{Cos}(2\pi ik/N) \quad i = 0\dots1023$$

$$\text{Sig4}(i) = \sum_{k=0}^{40} 10^{-51/40} \text{Cos}(20\pi ik/N) \quad i = 0\dots1023$$

Les figures (5.3) à (5.6) montrent les spectres de ces signaux : le spectre de Sig1 contient une raie, ceux de Sig2 et Sig3 contiennent 40 raies et celui de Sig4 contient 13 raies (au lieu de 40 raies)

Ces spectres à part le dernier sur lequel on reviendra, sont conformes à la théorie. L'évaluation de la qualité de ces spectres se fait en évaluant le bruit de calcul ou le rapport signal sur bruit. Pour ce faire, on a noté pour chaque étage si un débordement a eu lieu, ceci est donné par la table 5.1. L'évaluation du bruit se fait comme suit :

- En étant à l'étage m , m variant de 1 à $\text{Log}_2 N$, le bruit introduit a une variance de :

$$\sigma_{B1}^2 = \frac{2^{-2b}}{6}$$

s'il n'y a pas eu de débordement et une variance de :

$$\sigma_{B2}^2 = 7 \cdot \frac{2^{-2b}}{24}$$

si un débordement a eu lieu. b est le nombre de bits de la partie fractionnaire. Si un débordement a eu lieu dans un des deux premiers étages, sa variance sera :

$$\sigma_{B3}^2 = \frac{2^{-2b}}{4}$$

- Un bruit généré à l'étage m se trouve propagé vers la sortie par une constante. $N/2^r$, r est le nombre d'étages supérieur à m dans lesquels un débordement survient. Sachant que chaque noeud de sortie est connecté à $(N / 2^m)$ sources de bruit généré à l'étage m, la variance de ce bruit sera multipliée par $(N / 2^m) (\frac{1}{4})^r$. Sur cette base et en utilisant la table 5-1, on évalue le bruit de chaque signal. Les résultats sont reportés sur la table 5.2. La table contient aussi d'autres informations.

- Temps = temps de calcul à base d'un temps de cycle de 125 ns
- σ_s^2 = Puissance du signal à la sortie.
- nb-bits = nombre de bits affectés par le bruit.
- S/B = rapport signal sur bruit σ_s^2 / σ_B^2
- σ_R^2 = Puissance du bruit à la sortie.
- nb cycles = nombre de cycles processeur.

etage Signal	1	2	3	4	5	6	7	8	9	10
Sig1	1	0	1	1	1	1	1	1	1	1
Sig2	1	0	0	0	0	0	1	1	1	1
Sig3	0	0	0	0	0	1	1	1	1	1
Sig4	0	1	0	1	1	1	1	1	1	1

Table 5.1 : nombre de débordement par étage pour les signaux modèles considérés

	nb-cycle	Temps (ms)	$\sigma_B^2 (2^{2b})$	nb-bit	S/B db	σ_s^2
Sig1	96229	2.03	0.583	1	131.11	$1.953 \cdot 10^{-3}$
Sig2	79793	9.97	1.203	1	117.69	0.013274
Sig3	79793	9.97	0.601	1	123.71	0.013274
Sig4	92120	11.52	0.584	1	114.87	$3.0147 \cdot 10^{-4}$

Table 2 : Performances des signaux modeles

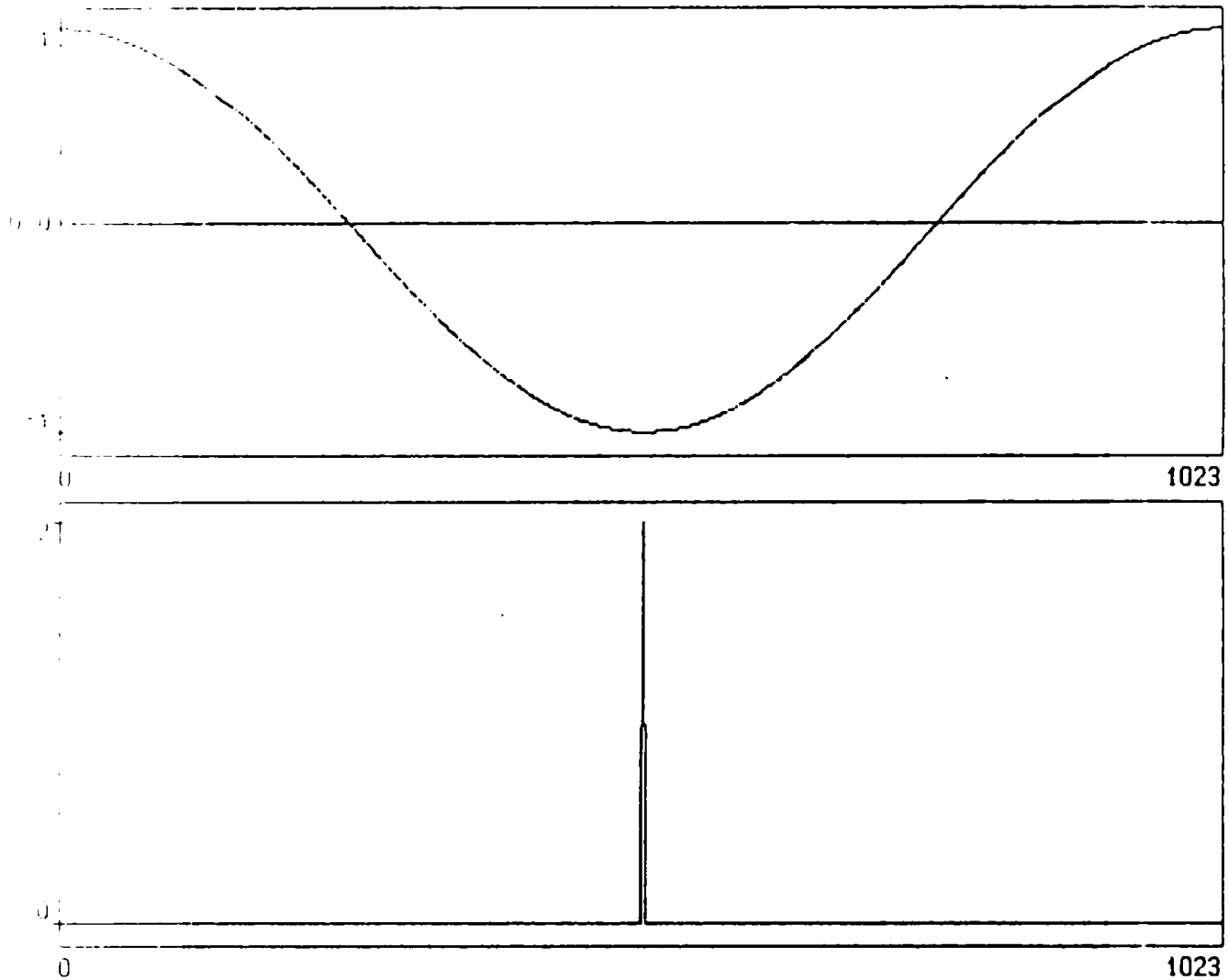


Figure 5.3 a) signal SIG1 b) spectre de SIG1 obtenu par CBFP

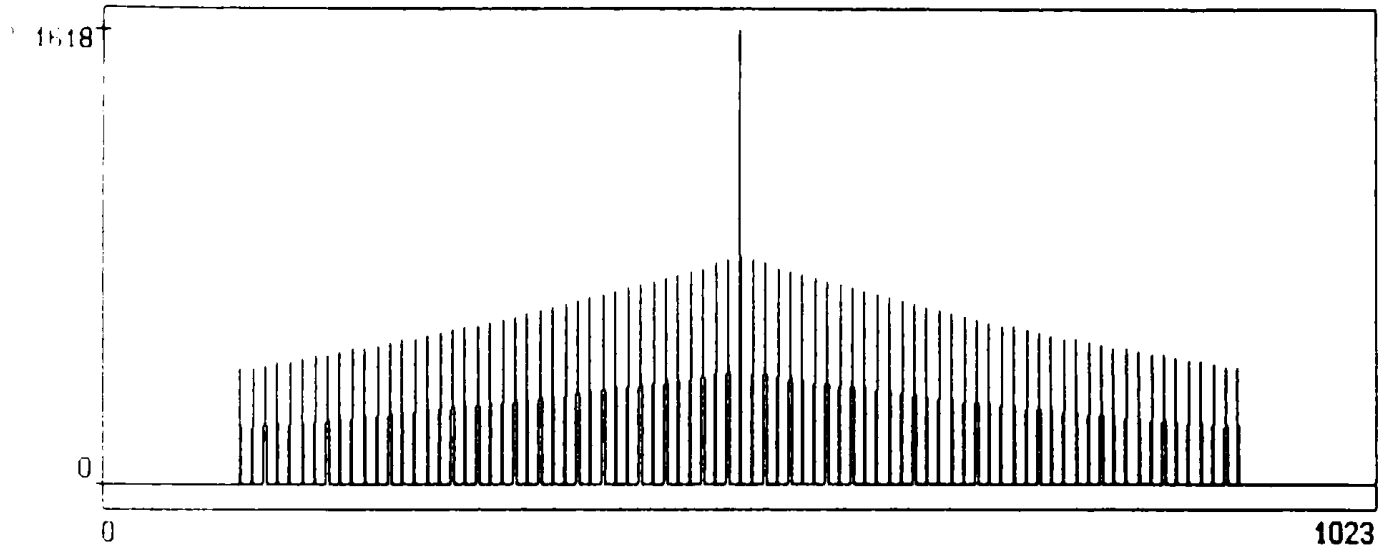
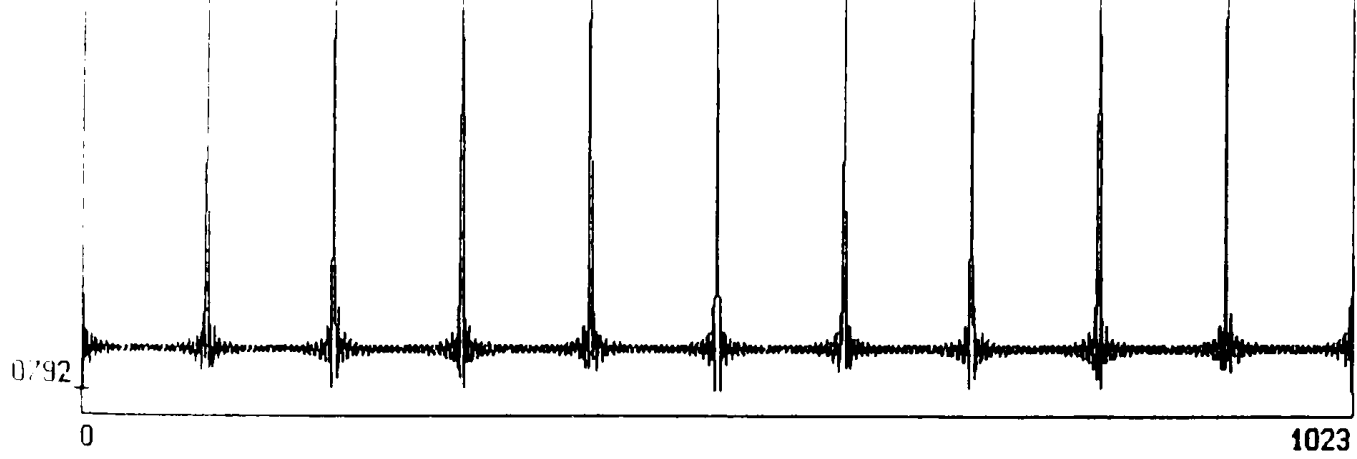


Figure 5.4 : a) signal SIG2 b) spectre de SIG2 obtenu par CBFP

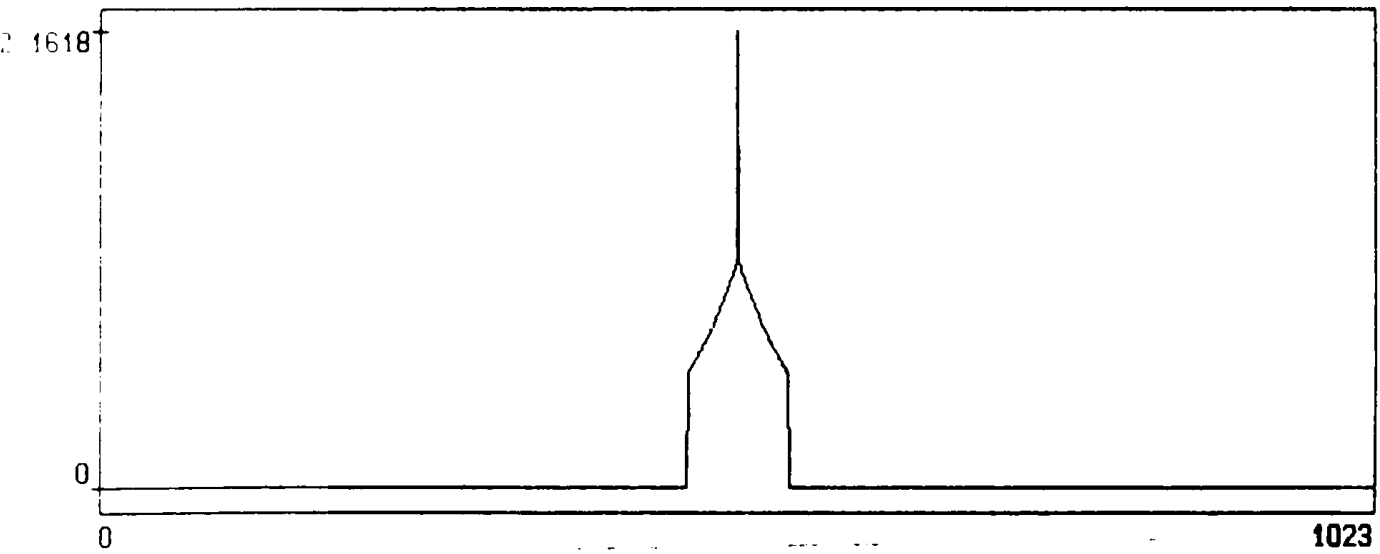
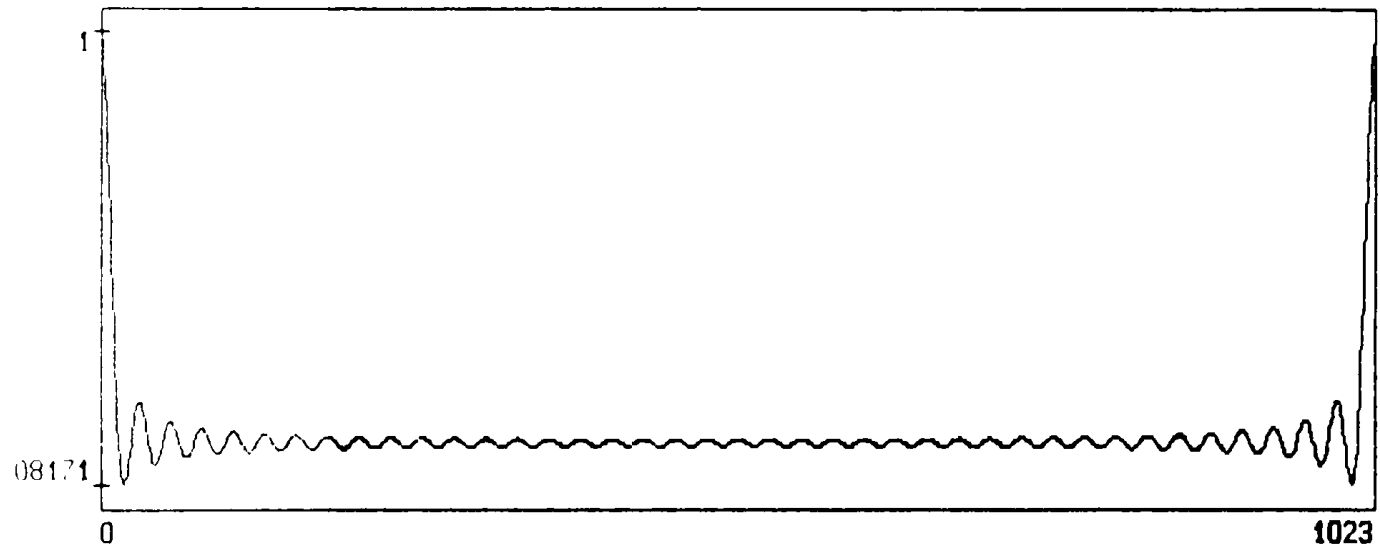


Figure 5.5 : a) signal SIG3 b) spectre de SIG3 obtenu par CBFP

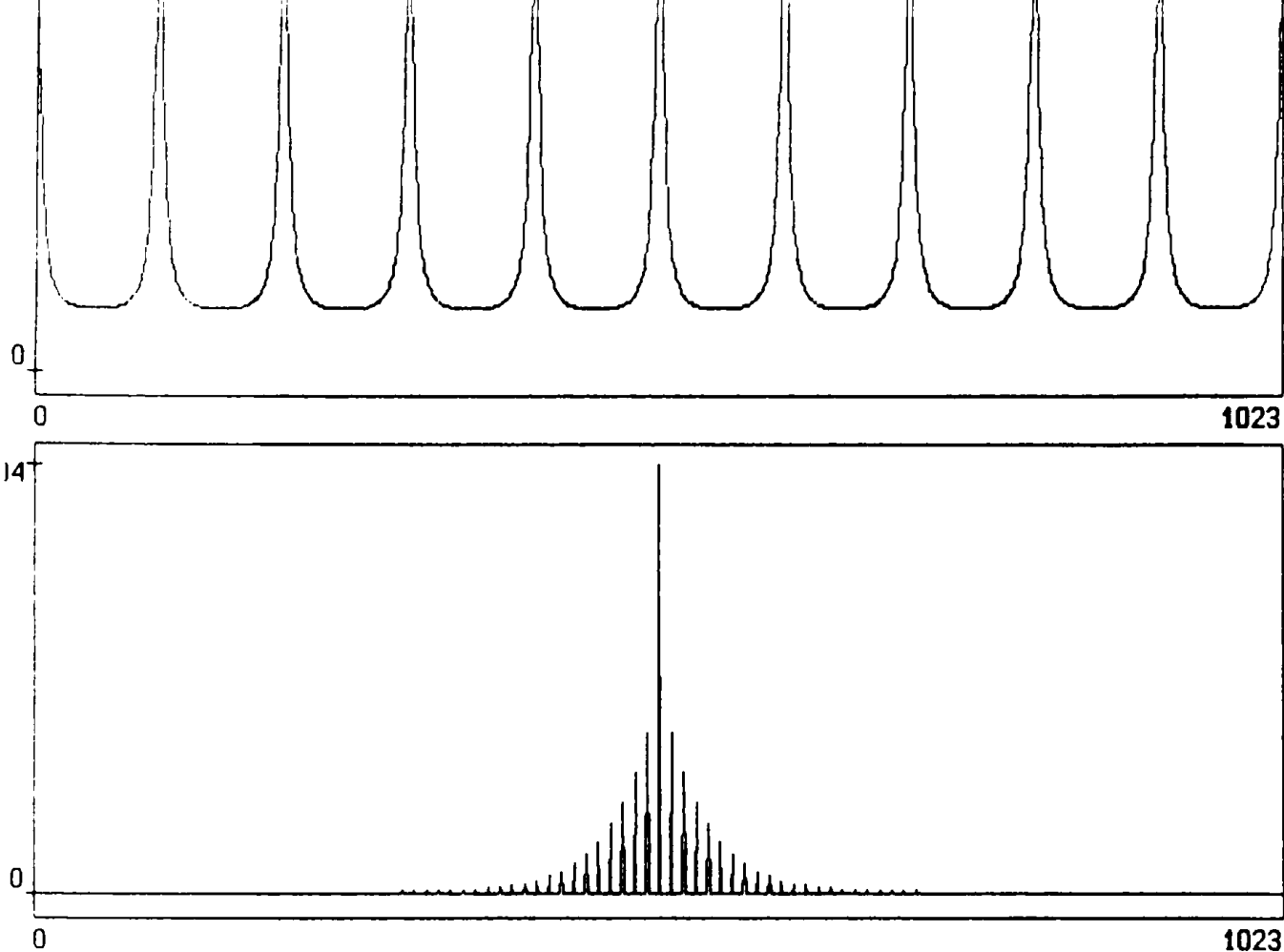


Figure 5.6 : a) signal SIG4 b) spectre de SIG4 obtenu par CBFP

5.2.2 La mise inconditionnelle en format flottant par bloc :

Dans cette méthode, les données réelles et imaginaires sont divisées par 2 à la sortie de chaque étage. Le programme correspondant se trouve dans le paragraphe A2.1.2 de l'annexe 2.

La méthode de programmation adoptée est la même que celle de la méthode précédente, i.e le premier et le deuxième étage sont calculés à part selon les équations (5.5) à (5.12).

Dynamique du signal d'entrée

Les papillons des deux premiers étages se calculent en huit cycles et ceux des autres étages se calculent en treize cycles.

Pour les même raisons que celles de la première méthode, la dynamique du signal se trouve réduite à 13 bits.

Temps de calcul

Le temps de calcul d'une FFT sur N points est le même pour n'importe quel signal. Il est donné par la formule :

$$T_{\text{exec}} = (6.5N \cdot \text{Log}_2 N + N + 19\text{Log}_2 N + 23)T_{\text{cycle}}$$

Soit pour une FFT sur 1024 points et un temps de cycle de 125 ns.

$$T_{\text{exec}} = 8,47 \text{ ms}$$

Bruit de calcul

Le bruit de calcul introduit par cette méthode de mise à l'échelle a une variance de :

$$\sigma_B^2 = \frac{1}{3} 2^{-2b}$$

Soit en nombre de bit [6]:

$$\text{nb-bit} = \text{INT} \left[\frac{1}{2} \text{Log}_2 (2^{2b} \sigma_B^2) \right] + 1$$

INT (.) : dénote la partie entière.

Soit nb-bit = 1 bit

Donc pour tout signal, on a toujours un bit qui est affecté par le bruit (le bit LSB).

Applications

On a appliqué le programme résultant de cette méthode aux signaux donnés au paragraphe précédent. Les spectres correspondants sont donnés par les figures (5.7.a) à (5.10.a). Les performances sont notées sur la table 5.3.

	nb-cycles	Temps ms	σ_s^2	$\sigma_B^2 (2^{-2b})$	nb-bit	S/B db
Sig1	67797	8.47	$4.8828 \cdot 10^{-4}$	1/3	1	123.93
Sig2	67797	8.47	$1.2962 \cdot 10^{-5}$	1/3	1	92.41
Sig3	67797	8.47	$1.2962 \cdot 10^{-5}$	1/3	1	92.41
Sig4	67797	8.47	$1.8842 \cdot 10^{-5}$	1/3	1	95.66

Table 5.3 : Performances des signaux modèles obtenus avec la mise à l'échelle inconditionnelle

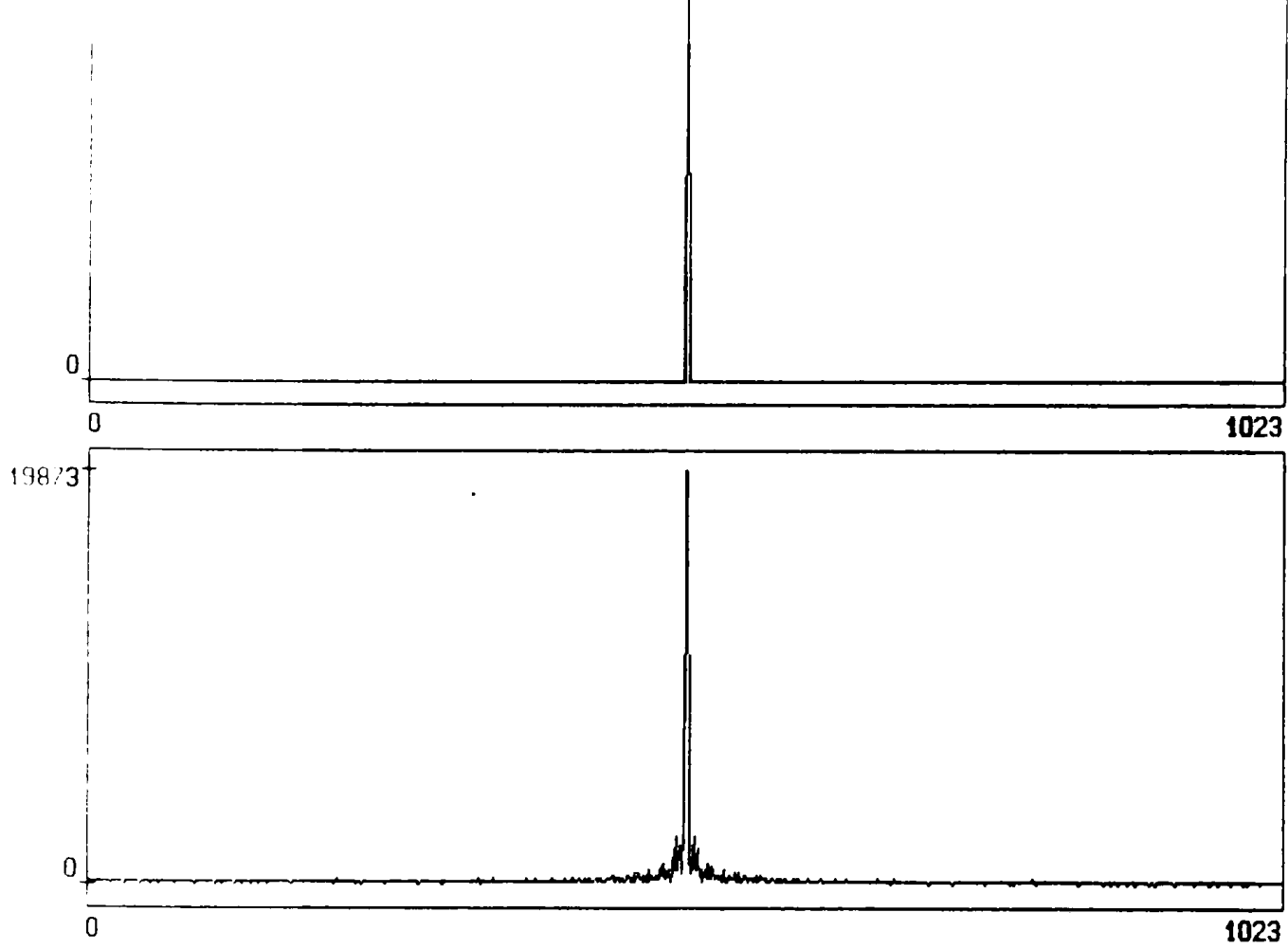


Figure 5.7 : spectre du signal SIG1 obtenu par a) UBFP b)INPUT

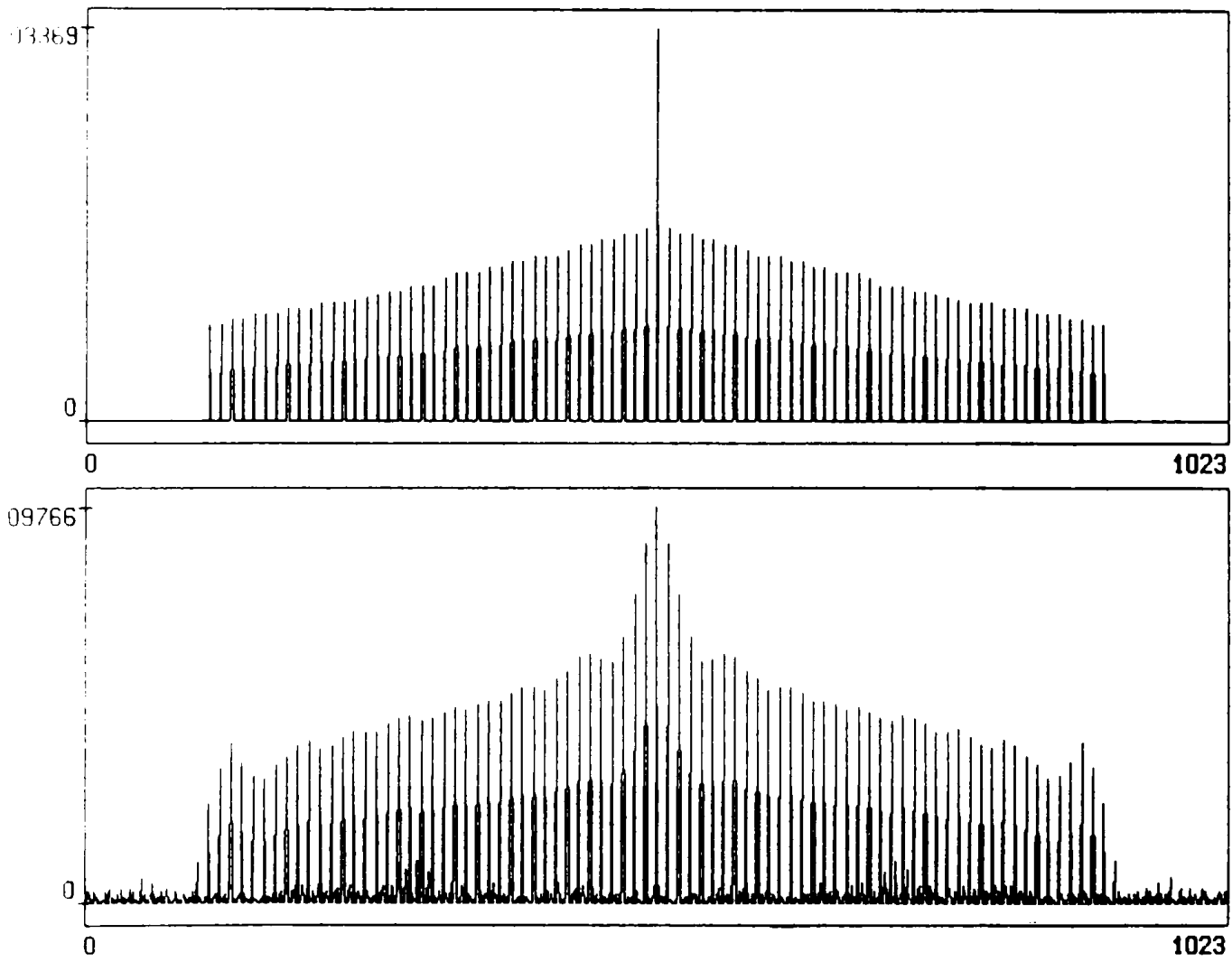


Figure 5.8 : spectre du signal SIG2 obtenu par a) UBFP b)INPUT

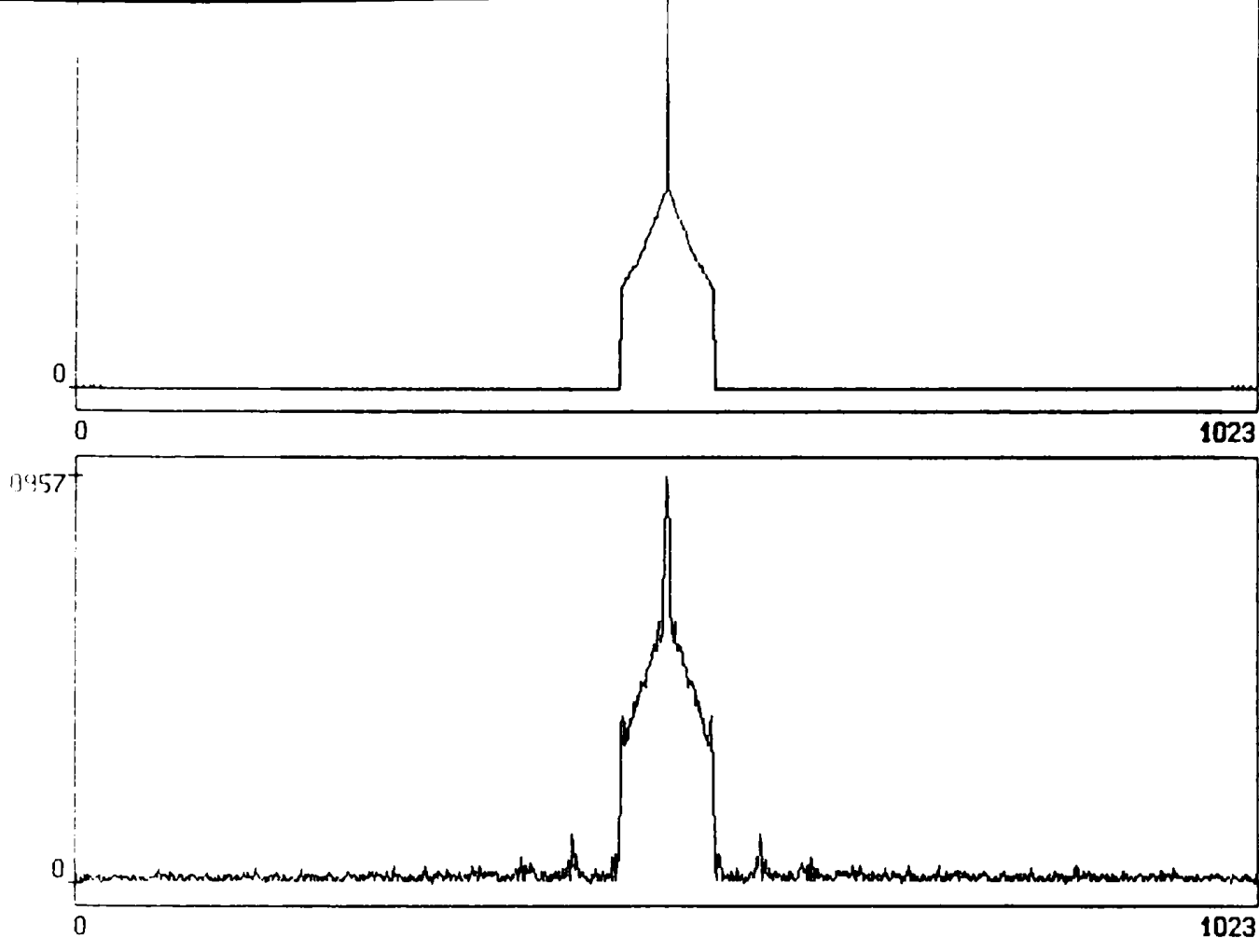


Figure 5.9 : spectre du signal SIG9 obtenu par a) UBFP b)INPUT

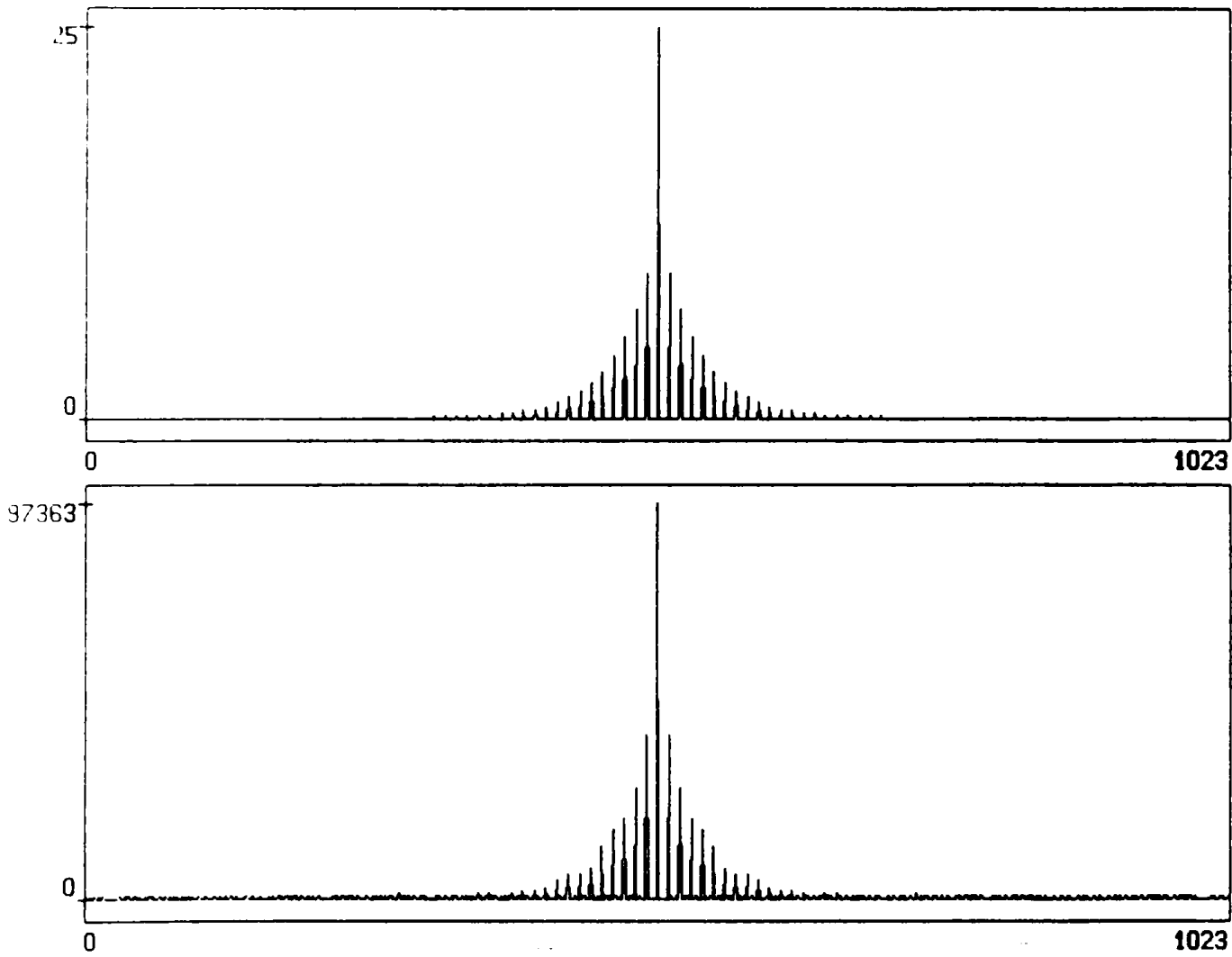


Figure 5.10 : spectre du signal SIG9 obtenu par a) UBFP b)INPUT

5.2.3. Mise à l'échelle des données à l'entrée

La troisième méthode de prévention contre le débordement consiste à diviser les parties réelles et imaginaires du signal par N , N étant le nombre d'échantillons. Le programme correspondant à cette méthode se trouve dans le paragraphe A2.1.3 de l'annexe 2.

Dynamique du signal d'entrée

Pour cette méthode, la dynamique du signal dépend du nombre d'échantillons N . Elle est de $(15 - \text{Log}_2 N)$ bits. Et plus le nombre d'échantillons augmente, plus la dynamique diminue et plus les résultats seront imprécis.

Temps de calcul

Les papillons des deux premiers étages se calculent en huit cycles et ceux des autres étages se calculent en neuf cycles.

Le temps de calcul d'une FFT sur N points ne dépend pas de la nature du signal. Ce temps est donné par la formule :

$$T_{\text{exec}} = (4.5 N \text{Log}_2 N + 5N + 19\text{Log}_2 N + 46) \cdot T_{\text{cycle}}$$

Pour un temps de cycle de 125 ns et 1024 points, on aura :

$$T_{\text{exec}} = 6.43 \text{ ms}$$

Bruit de calcul

Le bruit de calcul introduit par papillon est :

$$\sigma_{\text{Bp}}^2 = 2 \cdot \frac{2^{-2b}}{12}$$

Comme les deux premiers étages sont réalisés sans bruit, on aura $(\frac{N}{4} - 1)$ sources de bruit qui contribuent dans le bruit présent à chaque noeud de sortie de la FFT. Ceci donne un bruit total :

$$\sigma_{\text{B}}^2 = \frac{N}{24} 2^{-2b}$$

Soit, en nombre de bits [6] :

$$\text{nb-bit} = \text{INT} \left[\frac{1}{2} \text{Log}_2 (2^b \sigma_{\text{B}}^2) \right] + 1$$

$$\text{nb-bit} = \text{INT} \left[\frac{1}{2} \text{Log}_2 N - 2.29 \right] + 1$$

Pour une FFT, sur 1024 points, on aura :

$$\text{nb-bit} = 3 \text{ bits}$$

Applications

On a appliqué le programme résultant de cette méthode aux signaux précédents. Les spectres correspondants sont donnés aux figures (5.7.b) à (5.10.b). Les performances sont données sur la table 5.4.

	nb-cycles	Temps ms	σ_s^2	$\sigma_B^2(2^{2b})$	nb-bit	S/B db
Sig1	51436	6.43	$4.8828125 \cdot 10^{-4}$	42.667	3	81.79
Sig2	51436	6.43	$1.2962913 \cdot 10^{-5}$	42.667	3	50.27
Sig3	51436	6.43	$1.2962913 \cdot 10^{-5}$	42.667	3	50.27
Sig4	51436	6.43	$1.8842411 \cdot 10^{-5}$	42.667	3	53.52

Table 5.4 : Performances des signaux modèles obtenues avec la mise à l'échelle des données à l'entrée

5.2.4. Comparaison des trois méthodes :

En comparant les figures (5.3) à (5.10) et les tables (5.2) à (5.4) donnant respectivement les spectres des signaux modèles obtenus avec les trois méthodes de mise à l'échelle et les performances des trois méthodes pour ces signaux, on tire ce qui suit :

- La mise conditionnelle en format flottant offre le meilleur rapport signal sur bruit et la meilleure dynamique du signal donc la meilleure précision. Elle a pour inconvénient d'être la plus lente. Son temps de calcul peut être le double de celui de la troisième méthode. Un autre inconvénient très gênant est que son temps de calcul dépend du signal. Si on a à concevoir une application en temps réel contenant plusieurs tâches parmi lesquelles la FFT et qu'on doit attribuer un temps de calcul à chaque tâche, on aura un problème avec la FFT. Si on lui attribue

le temps maximum de calcul, on risque d'être très sévère puisqu'il se peut que la FFT s'exécute dans un temps beaucoup plus petit et on aura ainsi un temps énorme gâché alors qu'on pouvait l'exploiter à faire une autre tâche. Si on lui attribue un temps plus petit, on risque de tomber dans la situation inverse i.e la FFT va consommer un temps attribué à une autre tâche et ceci peut erroner les résultats.

Le choix du temps à attribuer à la FFT dépendra donc des statistiques du signal à traiter.

Cette méthode de mise à l'échelle doit être utilisée pour les applications où la précision prime sur le temps de calcul.

- La mise à l'échelle des données à l'entrée offre le temps de calcul le plus petit. Ce temps est connu avec précision. Son inconvénient est qu'elle présente le rapport signal sur bruit le plus défavorable (le plus faible). Les figures (5.7.b) à (5.10.b) montrent clairement l'effet du bruit de calcul sur les spectres des signaux.

Cette méthode doit être utilisée pour les applications où le temps de calcul prime sur la précision.

- La mise inconditionnelle en format flottant offre une solution intermédiaire entre les deux premières. Elle est plus rapide et moins précise que la première et plus précise et moins rapide que la seconde. Elle sera la méthode la plus utilisée si on n'a pas une contrainte sévère sur le temps de calcul ou sur la précision.

Un dernier point à noter est la limite de la représentation en virgule fixe. Les figures (5.6) et (5.10) donnent les spectres du signal Sig4. Celui-ci contient en théorie 40 raies alors que ces figures montrent moins de 40 raies. Ceci est une conséquence de la représentation en virgule fixe. Si b est le nombre de bits de la partie fractionnaire, les harmoniques dont l'amplitude est inférieure à 2^{-b} n'apparaîtront pas dans le résultat final. Pour les applications nécessitant une très grande précision de calcul, la virgule fixe est très limitée. Cette précision sera atteinte par une représentation en multiprécision ou par une représentation en virgule flottante.

5.2.5 La FFT d'une séquence réelle

On sait d'après le paragraphe 2.8 que la FFT d'une séquence réelle de durée N peut être déduite de celle d'une séquence complexe de durée $N/2$. Ceci fait gagner près de la moitié du temps de calcul total.

Le calcul se fait selon les trois étapes suivantes :

1- Dispersion des données : construction du signal complexe de durée $N/2$ selon la relation (2.24)

2- FFT du signal complexe : calcul de la FFT du signal complexe de durée $N/2$.

3- Rassemblement des données : déduction de la TFD du signal réel à partir de celle du signal complexe selon la relation (2.30).

Soit $x(n)$, le signal réel et $X(k)$ sa TFD et soit $z(n)$ le signal complexe et $Z(k)$ sa TFD. Le rassemblement des données se fait selon les équations suivantes :

$$X_R(k) = \frac{1}{2} \{ [Z_R(k) + Z_R(\frac{N}{2} - k)] - \sin \frac{2\pi k}{N} [Z_R(k) - Z_R(\frac{N}{2} - k)] \\ + \cos \frac{2\pi k}{N} [Z_I(k) + Z_I(\frac{N}{2} - k)] \}$$

$$X_I(k) = \frac{1}{2} \{ [Z_I(k) - Z_I(\frac{N}{2} - k)] - \cos \frac{2\pi k}{N} [Z_R(k) - Z_R(\frac{N}{2} - k)] \\ - \sin \frac{2\pi k}{N} [Z_I(k) + Z_I(\frac{N}{2} - k)] \}$$

$$X_R(N-k) = X_R(k) \quad k = 0, \dots, \frac{N}{2} - 1 \quad (5.13)$$

$$X_I(N-k) = -X_I(k)$$

Les indices I et R dénotent les parties imaginaire et réelle.

Les deux dernières équations de la relation (5.13) sont dues à la propriété de parité de la TFD d'une séquence réelle (paragraphe 2.4.3)

Les programmes de dispersion et de rassemblement des données sont donnés au paragraphe A2.1.4 de l'annexe 2.

La table 5.5 donne une comparaison entre les FFT complexe et réelle des signaux modèles Sig1, Sig2, Sig3 et Sig4, selon les trois méthodes de mise à l'échelle.

Cette table met en évidence l'avantage d'exploiter le fait que les signaux sont réels. Car ce sont ces signaux qu'on a, la plus part du temps, à traiter.

	CBFP		UBFP		INPUT	
	Réelle	Complexe	Réelle	Complexe	Réelle	Complexe
Sig 1	6.11	12.03	4.80	8.47	3.65	6.43
Sig 2	5.08	9.97	4.80	8.47	3.65	6.43
Sig 3	5.08	9.97	4.80	8.47	3.65	6.43
Sig 4	5.89	11.52	4.80	8.47	3.65	6.43

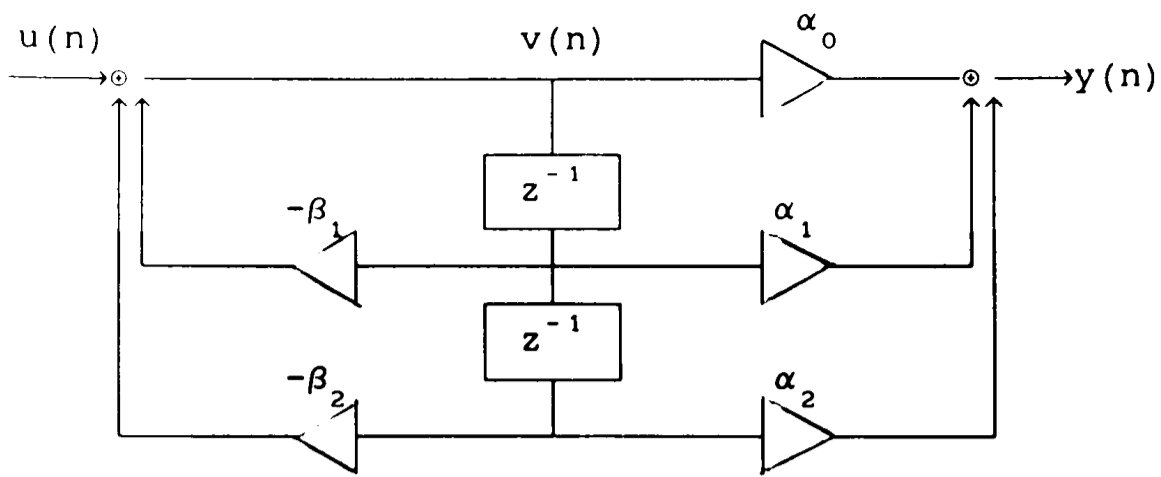
Table 5.5 : Comparaison entre les temps de calcul des FFT complexe et réelle des signaux modèles (en ms)

5.2.6 Opération d'inversion de bits

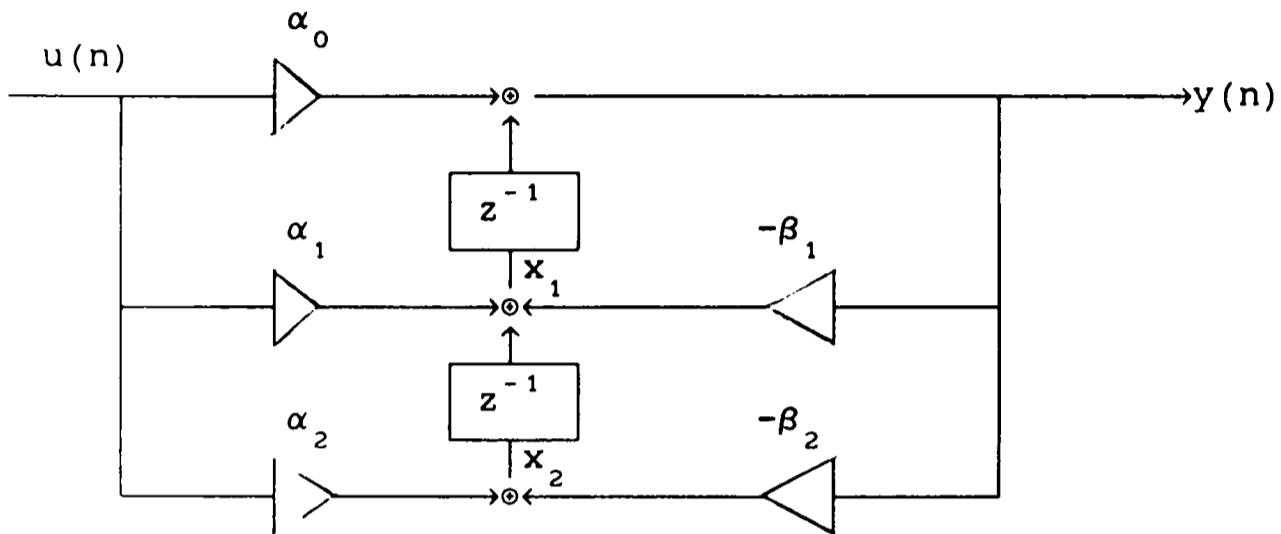
Avant de calculer la FFT, les séquences doivent être placées dans l'ordre bit-inversé. L'ADSP-2100 possède une option d'inversion des lignes d'adresse dédiée au calcul de la FFT. Cette option évite l'inversion des bits d'adresse par le calcul et fait gagner un temps de calcul très important. Le programme est donné dans le paragraphe A2.1.5 de l'annexe2.

5.3. Le filtrage IIR

On a montré dans le troisième chapitre que la structure cascade des filtres IIR est la moins sensible aux erreurs de quantification des coefficients. Et, on a noté, que les cellules du deuxième ordre les plus utilisées sont les formes 1D et 2D. Dans le quatrième chapitre, on a mis à l'échelle ces deux formes de cellules pour les prévenir contre les débordements qui peuvent survenir au cours des calculs. Les formes dérivées correspondent à la figure 5.11. Les programmes en langage assembleur de l'ADSP-2100 réalisant le filtrage selon les deux formes de cellules se trouvent dans le paragraphe A2.2 de l'annexe 2.



a) Forme 1D



b) Forme 2D

Figure 5.11 : Cellule cascade mise a l'echelle

5.3.1 Structure 1D :

De la figure 5.11.a, on tire différentes variables de la cellule : $u(n)$ et $y(n)$ sont respectivement l'entrée et la sortie de la cellule, X_1 et X_2 sont les états des calculs précédents, $v(n)$ est une variable intermédiaire et α_0 , α_1 , α_2 , β_1 et β_2 sont les coefficients de la cellule.

Les équations liant ces coefficients variables sont :

$$\begin{aligned}
 v(n) &= u(n) - \beta_1 \cdot X_1 - \beta_2 \cdot X_2 \\
 y(n) &= \alpha_0 \cdot v(n) + \alpha_1 \cdot X_1 + \alpha_2 \cdot X_2 \\
 X_2 &= X_1 \\
 X_1 &= v(n)
 \end{aligned}
 \tag{5.14}$$

Le calcul d'une cellule se fait selon les étapes suivantes :

- Calculer $v(n)$
- Calculer $y(n)$
- Mettre à jour X_2
- Mettre à jour X_1

Le point important qui reste à définir est le format des différentes variables :

Les coefficients : On a généralement :

$$|\beta_1| \leq 2 \quad \text{et} \quad |\alpha_1| \leq 1$$

donc 2 bits seront nécessaires pour la partie entière et 14 bits pour la partie fractionnaire pour une longueur de mot de 16 bits. Ce format sera noté format (2.14).

Les données : $u(n)$, X_1 et X_2 sont supposées normalisées entre -1 et +1. Ceci nous conduit à :

$$|v(n)| \leq 5$$

Donc, trois bits seront nécessaires pour représenter la partie entière des données. En plus du bit signe, ceci nous conduit à ce que les données sont représentées dans le format (4.12).

Des conversions de format seront nécessaires lors des calculs intermédiaires des équations de la relation (5.14).

Quand on multiplie une donnée sous le format (2.14) par une donnée sous le format (4.12), on aura un résultat sous le format (6.10). L'ADSP-2100 dispose d'un multiplieur qui fait un décalage automatique de un bit vers la gauche après chaque multiplication. Soit donnant un résultat sous le format (5.11). C'est sous ce format, que sera le résultat de $(\beta_1 \cdot X_1 + \beta_2 \cdot X_2)$. Pour pouvoir l'additionner à $u(n)$ qui est sous le format (4.12), on doit le décaler de 1 bit vers la gauche pour rectifier la position de la virgule.

La même chose est notée concernant la sortie $y(n)$ de la cellule qui doit être décalée de 1 bit vers la gauche pour pouvoir être

l'entree de la prochaine cellule.

Le programme développé tient compte des équations de la relation (5.14) et des conversions de format notées ci-dessus. Une cellule du second ordre est calculée en 9 cycles processeur et contient 2 sources de bruit.

5.3.2 Structure 2D :

De la figure 5.10.b, on tire les variables de la cellule : $u(n)$ et $y(n)$ sont respectivement l'entrée et la sortie de la cellule, X_1 et X_2 sont les états des calculs précédents, α_0 , α_1 , α_2 , β_1 et β_2 sont les coefficients de la cellule.

Les équations liant ces différentes variables sont :

$$\begin{aligned}y(n) &= \alpha_0 \cdot u(n) + X_1 \\X_1 &= \alpha_1 \cdot u(n) - \beta_1 \cdot y(n) + X_2 \\X_2 &= \alpha_2 \cdot u(n) - \beta_2 \cdot y(n)\end{aligned}\tag{5.14}$$

Le calcul de la cellule se fait selon les étapes suivantes :

- Calculer $y(n)$
- Calculer X_1
- Mettre à jour X_1
- Mettre à jour X_2

Le format des différentes variables se déduit de la même façon que pour la structure 1D. On trouve que les coefficients sont dans le format (2.14) et les données dans le format (4.12). Des conversions de format sont nécessaires après les multiplications pour pouvoir faire les différentes additions. Ceci est fait par des décalages.

Le programme correspondant à une cellule s'exécute en 13 cycles et contient 3 sources de bruit.

5.3.3 Comparaison des deux structures :

La cellule réalisée en structure 1D s'exécute en 9 cycles et contient 2 sources de bruit. Par contre, la structure 2D s'exécute en 13 cycles et contient 3 sources de bruit.

Donc la structure 1D est la meilleure à utiliser avec

5.3.4 Temps de calcul :

Si N est l'ordre du filtre, un échantillon du signal de sortie est calculé en $(4.5 N + 13)$ cycles. $(N + 16)$ cycles sont nécessaires pour initialiser le programme.

Donc, la période d'échantillonnage minimale sera $(5.5N+29)T_c$, T_c étant le temps de cycle du processeur. Si on a $N=10$ et $T_c=125\text{ns}$, la fréquence d'échantillonnage maximale sera de 95 K Hz. La relation entre la fréquence d'échantillonnage F_e et l'ordre du filtre est :

$$F_e \leq 1/(5.5 N + 29)T_c$$

5.3.5. Applications :

On a considéré 3 applications : Les deux premières étant le filtrage du signal modèle Sig2 décrit dans le paragraphe 5.2, une fois par un filtre passe-bas et une fois par un filtre passe-haut. La troisième application est le filtrage d'un signal EEG contenant un bruit additif. Les coefficients des filtres ont été obtenus par le logiciel MONARCH [39]. Le signal EEG bruité est donné dans la disquette de distribution du même logiciel. Les filtres sont synthétisés selon l'approximation de Butterworth.

Le logiciel MONARCH est destiné pour les chercheurs dans le domaine du traitement de signal. Il permet, entre autre, la synthèse des filtres IIR selon les approximations de Butterworth, de Tchebyshev et de Cauey, la synthèse des filtres FIR selon l'algorithme de Remez ainsi que d'autres opérations de traitement numérique du signal.

a) Filtre passe-bas

La figure (5.12) donne les spécifications, l'ordre et les coefficients du filtre. Les figures (5.13.a) et (5.13.b) donnent sa réponse impulsionnelle et sa réponse fréquentielle. Les figures (5.13.c) et (5.13.d) donnent le spectre original de Sig2 et le spectre du signal filtré. On remarque que les raies du spectre de Sig2 qui sont en dehors de la bande passante du filtre ne se trouve plus dans le spectre du signal.

b) Filtre passe-haut

La figure (5.14) donne les spécifications, l'ordre et les coefficients du filtre. Les figures (5.15.a), (5.15.b), (5.15.c) et (5.15.d) donnent respectivement la réponse impulsionnelle du filtre sa réponse fréquentielle, le spectre original de Sig2 et son spectre filtré.

c) Signal EEG

Ce signal est prélevé à une fréquence de 166 Hz. Il contient un bruit additif de 60 Hz dû au secteur (norme américaine). On le filtre par un filtre passe-bas ayant une bande passante entre 0 et 55 Hz et une bande de transition entre 55 et 60 Hz.

L'atténuation dans la bande passante est de 3 db et dans la bande atténuée de 20 db. Le filtre résultant est du 10^{ème} ordre. Les coefficients du filtre résultant sont donnés dans la figure 5.16. Ses réponses impulsionnelle et fréquentielle sont données dans la figure 5.17. Le signal EEG bruité, son spectre, le signal filtré et le spectre du signal filtré sont donnés dans la figure 5.18.

La figure 5.18.a montre l'effet du bruit additif sur le signal. La figure 5.18.b montre un pic situé à une fréquence de 60 Hz dû au bruit. La figure 5.18.c montre le signal EEG filtré et la figure 5.18.d montre son spectre. On voit que le pic dû au bruit additif n'existe plus dans le spectre du signal filtré.

Bruit de calcul

Le calcul du bruit à la sortie des différents filtres a été omis. On ne dispose pas actuellement de logiciel qui optimise l'organisation des cellules du second ordre et calcule le bruit de quantification résultant.

```

Approximation      : Butterworth
Filter Function    : lp
PassBand Atten. AP (dB) :      3.00
StopBand Atten. AA (dB) :     20.00
Sampling Frequency :     166.00
PassBand Edge (Hz) :     55.00
StopBand Edge (Hz) :     60.00

```

DEVELOPED DIGITAL FILTER
Butterworth LowPass FILTER

ORDER : 10

$$H(Z) = H * H_1(Z) * H_2(Z) * \dots$$

$$H\#(Z) = (A0\# + A1\#*Z + A2\#*Z^2) / (B0\# + B1\#*Z + B2\#*Z^2)$$

H = 2.697919959389123e-002

SECTION : 1 H1(Z)

A 0 1	:	1.0000000000000000e+000	B 0 1	:	7.444527897352987e-002
A 1 1	:	2.0000000000000000e+000	B 1 1	:	5.256343515026775e-001
A 2 1	:	1.0000000000000000e+000	B 2 1	:	1.0000000000000000e+000

SECTION : 2 H2(Z)

A 0 2	:	1.0000000000000000e+000	B 0 2	:	1.254269774550810e-001
A 1 2	:	2.0000000000000000e+000	B 1 2	:	5.505753443519887e-001
A 2 2	:	1.0000000000000000e+000	B 2 2	:	1.0000000000000000e+000

SECTION : 3 H3(Z)

A 0 3	:	1.0000000000000000e+000	B 0 3	:	2.370781972595001e-001
A 1 3	:	2.0000000000000000e+000	B 1 3	:	6.051967547345127e-001
A 2 3	:	1.0000000000000000e+000	B 2 3	:	1.0000000000000000e+000

SECTION : 4 H4(Z)

A 0 4	:	1.0000000000000000e+000	B 0 4	:	4.327120582635012e-001
A 1 4	:	2.0000000000000000e+000	B 1 4	:	7.009037020059861e-001
A 2 4	:	1.0000000000000000e+000	B 2 4	:	1.0000000000000000e+000

SECTION : 5 H5(Z)

A 0 5	:	1.0000000000000000e+000	B 0 5	:	7.598872811769465e-001
A 1 5	:	2.0000000000000000e+000	B 1 5	:	8.609626082056099e-001
A 2 5	:	1.0000000000000000e+000	B 2 5	:	1.0000000000000000e+000

figure 5.12: Gabarit du filtre passe-bas IIR

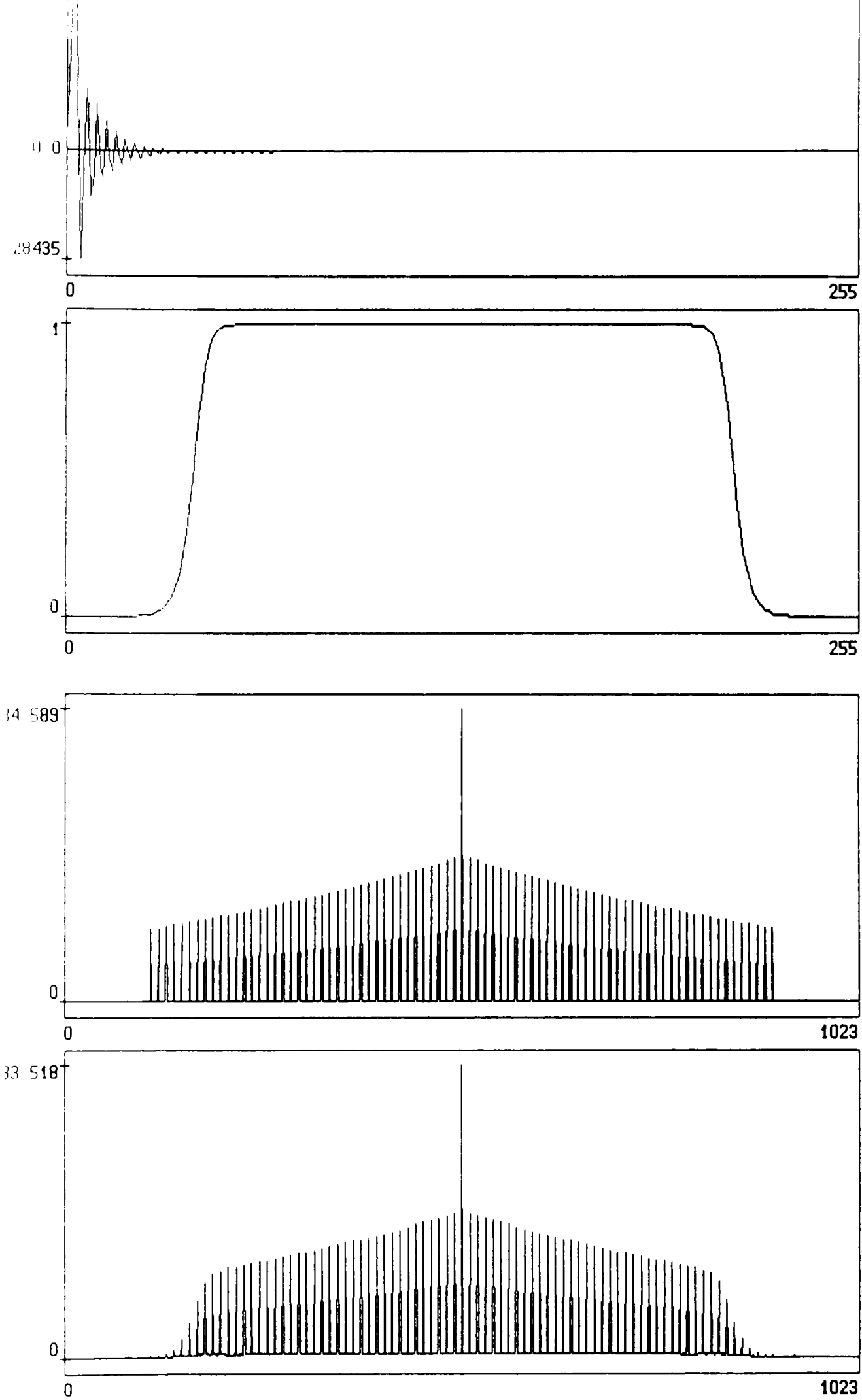


Figure 5.13 : a)réponse impulsionnelle du filtre passe-bas IIR
 b)réponse fréquentielle du filtre c)spectre du signal original
 d)spectre du signal filtré

```

Approximation      : Butterworth
Filter Function    : hp
PassBand Atten. AP (dB) :      3.00
StopBand Atten. AA (dB) :     20.00
Sampling Frequency :     160.00
PassBand Edge (Hz) :     40.00
StopBand Edge (Hz) :     35.00

```

DEVELOPED DIGITAL FILTER
Butterworth HighPass FILTER

ORDER : 12

$$H(Z) = H * H_1(Z) * H_2(Z) * \dots$$

$$H\#(Z) = (A_0\# + A_1\# * Z + A_2\# * Z^2) / (B_0\# + B_1\# * Z + B_2\# * Z^2)$$

H = 9.053750955351291e-004

SECTION : 1 H1(Z)

A 0 1	:	1.0000000000000000e+000	B 0 1	:	4.295955306027511e-003
A 1 1	:	-2.0000000000000000e+000	B 1 1	:	-1.987226925248741e-004
A 2 1	:	1.0000000000000000e+000	B 2 1	:	1.0000000000000000e+000

SECTION : 2 H2(Z)

A 0 2	:	1.0000000000000000e+000	B 0 2	:	3.956613966965216e-002
A 1 2	:	-2.0000000000000000e+000	B 1 2	:	-2.057016970360597e-004
A 2 2	:	1.0000000000000000e+000	B 2 2	:	1.0000000000000000e+000

SECTION : 3 H3(Z)

A 0 3	:	1.0000000000000000e+000	B 0 3	:	1.152292035189113e-001
A 1 3	:	-2.0000000000000000e+000	B 1 3	:	-2.206733472686136e-004
A 2 3	:	1.0000000000000000e+000	B 2 3	:	1.0000000000000000e+000

SECTION : 4 H4(Z)

A 0 4	:	1.0000000000000000e+000	B 0 4	:	2.431924204250813e-001
A 1 4	:	-2.0000000000000000e+000	B 1 4	:	-2.459937668853560e-004
A 2 4	:	1.0000000000000000e+000	B 2 4	:	1.0000000000000000e+000

SECTION : 5 H5(Z)

A 0 5	:	1.0000000000000000e+000	B 0 5	:	4.464627000089745e-001
A 1 5	:	-2.0000000000000000e+000	B 1 5	:	-2.862153938427372e-004
A 2 5	:	1.0000000000000000e+000	B 2 5	:	1.0000000000000000e+000

SECTION : 6 H6(Z)

A 0 6	:	1.0000000000000000e+000	B 0 6	:	7.690877206418854e-001
A 1 6	:	-2.0000000000000000e+000	B 1 6	:	-3.500540585684428e-004
A 2 6	:	1.0000000000000000e+000	B 2 6	:	1.0000000000000000e+000

Figure 5.14: Gabarit du filtre passe-haut IIR

```

Approximation : Butterworth
Filter Function : lp
PassBand Atten. AP (dB) : 3.00
StopBand Atten. AA (dB) : 20.00
Sampling Frequency : 166.00
PassBand Edge (Hz) : 55.00
StopBand Edge (Hz) : 60.00

```

DEVELOPED DIGITAL FILTER
Butterworth LowPass FILTER

ORDER : 10

$$H(Z) = H * H1(Z) * H2(Z) * \dots$$

$$H\#(Z) = (A0\# + A1\# * Z + A2\# * Z^2) / (B0\# + B1\# * Z + B2\# * Z^2)$$

H = 2.697919959389123e-002

SECTION : 1 H1(Z)

A 0 1	:	1.0000000000000000e+000	B 0 1	:	7.444527897352987e-002
A 1 1	:	2.0000000000000000e+000	B 1 1	:	5.256343515026775e-001
A 2 1	:	1.0000000000000000e+000	B 2 1	:	1.0000000000000000e+000

SECTION : 2 H2(Z)

A 0 2	:	1.0000000000000000e+000	B 0 2	:	1.254269774550810e-001
A 1 2	:	2.0000000000000000e+000	B 1 2	:	5.505753443519887e-001
A 2 2	:	1.0000000000000000e+000	B 2 2	:	1.0000000000000000e+000

SECTION : 3 H3(Z)

A 0 3	:	1.0000000000000000e+000	B 0 3	:	2.370781972595001e-001
A 1 3	:	2.0000000000000000e+000	B 1 3	:	6.051967547345127e-001
A 2 3	:	1.0000000000000000e+000	B 2 3	:	1.0000000000000000e+000

SECTION : 4 H4(Z)

A 0 4	:	1.0000000000000000e+000	B 0 4	:	4.327120582635012e-001
A 1 4	:	2.0000000000000000e+000	B 1 4	:	7.009037020059861e-001
A 2 4	:	1.0000000000000000e+000	B 2 4	:	1.0000000000000000e+000

SECTION : 5 H5(Z)

A 0 5	:	1.0000000000000000e+000	B 0 5	:	7.598872811769465e-001
A 1 5	:	2.0000000000000000e+000	B 1 5	:	8.609626082056099e-001
A 2 5	:	1.0000000000000000e+000	B 2 5	:	1.0000000000000000e+000

Figure 5.16: Gabarit du filtre passe-bas IIR pour le signal EEG

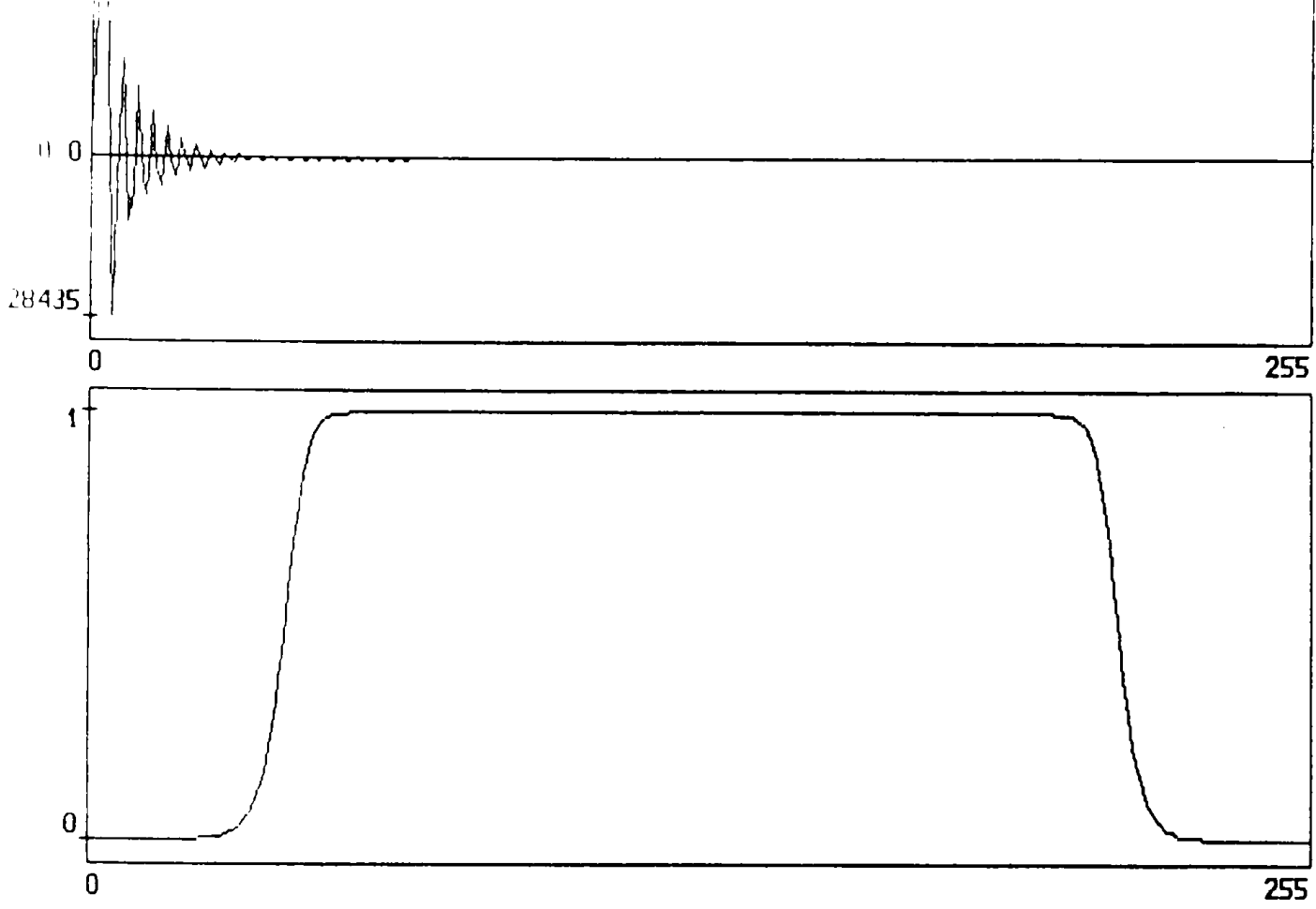


Figure 5.17 : a) réponse impulsionnelle du filtre passe-haut IIR pour le signal EEG b) réponse fréquentielle du filtre

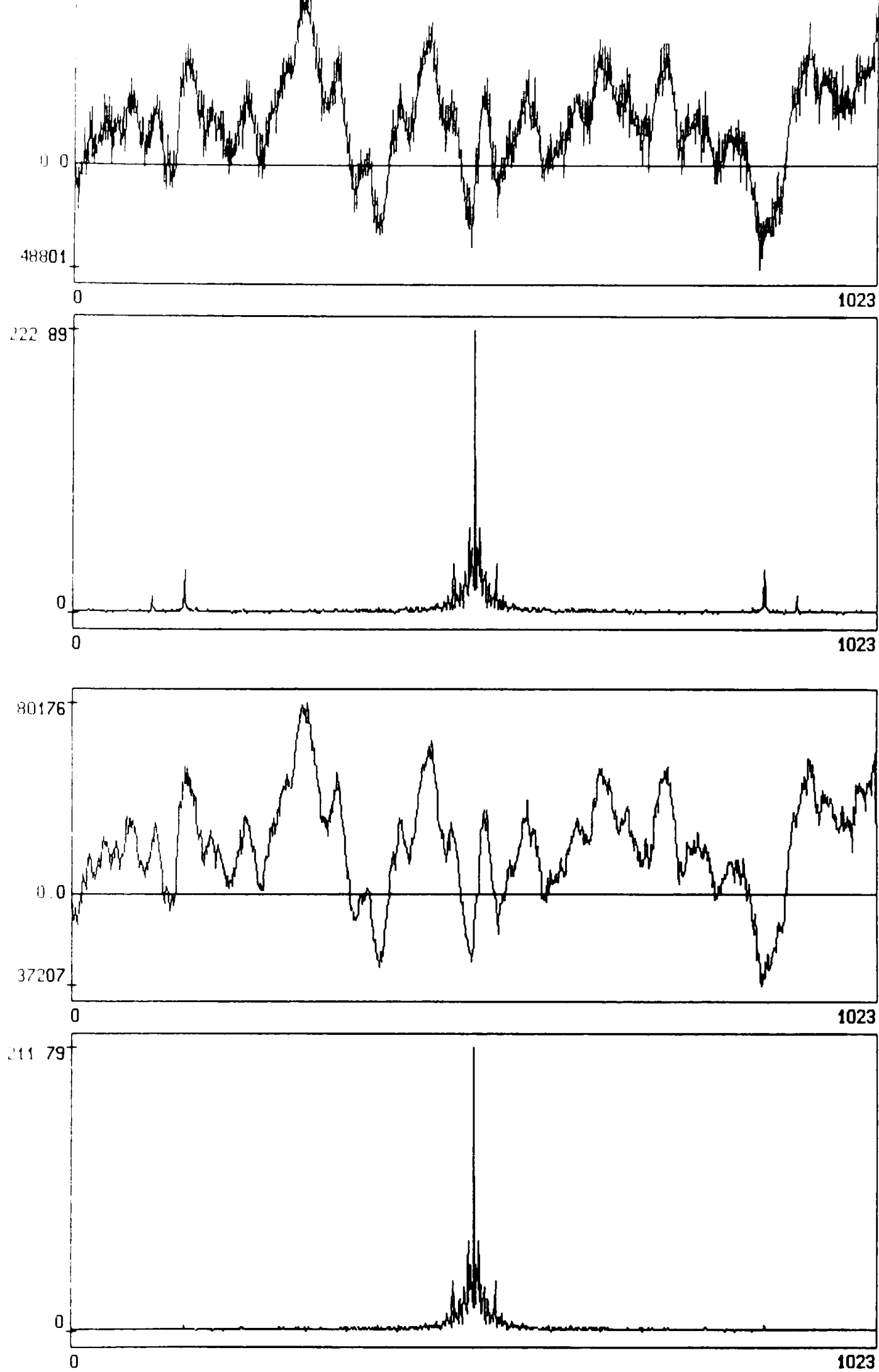


Figure 5.18 : a) signal EEG bruité b) spectre du signal EEG bruité
 c) signal EEG filtré par le filtre IIR d) spectre du EEG filtré

5.4 Le filtrage FIR :

La structure la plus utilisée pour les filtres FIR est la structure directe. Elle est basée sur le produit de convolution entre les coefficients du filtre et le signal d'entrée. Le programme réalisant le filtrage utilise un tableau contenant les coefficients du filtre et un tableau contenant des échantillons successifs du signal d'entrée qu'on appellera ligne à retard. A chaque fois qu'on fait l'acquisition d'un échantillon du signal, on le fait entrer dans la ligne à retard, on fait ensuite le produit cummulatif entre les éléments de la ligne à retard et les coefficients du filtre. A la fin des calculs, on arrondit le résultat et on met à jour la ligne à retard. Ceci rend négligeable le bruit de calcul présent à la sortie dont la variance est :

$$\sigma_B^2 = \frac{2^{-2b}}{12} = 7,76.10^{-11}$$

Ce qui fait qu'un seul bit sera affecté par le bruit.

Le programme correspondant est donné dans le paragraphe A2.3 de l'annexe 2.

5.4.1 Temps de calcul :

Si N est l'ordre du filtre, un échantillon de sortie est calculé en (N+13) cycles. (N+15) cycles sont nécessaires pour l'initialisation du programme.

Donc, la période d'échantillonnage minimale sera $(2N+28)T_c$, T_c étant le temps de cycle du processeur. Si, on a $N = 40$ et $T_c = 125$ ns, la fréquence d'échantillonnage maximale sera de 74 K Hz.

La relation entre la fréquence d'échantillonnage F_e et l'ordre du filtre est :

$$F_e \leq 1/(2N+28)T_c$$

5.4.2 Applications :

On a considéré 2 applications : le filtrage du signal Sig2 par un filtre passe-bas et le filtrage du signal EEG bruité. Les coefficients des filtres ont été obtenus par le logiciel MONARCH.

a) Filtre passe-bas

La figure (5.19) donne les coefficients du filtre, la figure

(5.20) donne les réponses impulsionnelle et fréquentielle du filtre. La figure (5.21) donne le spectre du signal Sig2 et le spectre du signal filtré.

b) Signal EEG

Le même signal EEG est filtré, cette fois par un filtre passe-bas FIR ayant une bande passante entre 0 et 55 Hz et une bande de transition entre 55 et 60 Hz.

La figure (5.22) donne les coefficients du filtre. La figure (5.23) donne les réponses impulsionnelle et fréquentielle du filtre. La figure (5.24) donne le signal bruité, son spectre, le signal filtré et son spectre. Les mêmes remarques que pour le filtrage IIR sont notées ici.

FINITE IMPULSE RESPONSE FOR A
LINEAR PHASE FILTER

BAND	LOWER BAND EDGE	UPPER BAND EDGE	DESIRED VALUE
1	0.000000e+000	3.000000e-001	1.000000e+000
2	3.300000e-001	5.000000e-001	0.000000e+000

FILTER LENGTH = 40

IMPULSE RESPONSE

```

H( 1)=  9.148428070751818e-003  =H(40)
H( 2)= -2.552111825093688e-002  =H(39)
H( 3)= -6.481305667170985e-003  =H(38)
H( 4)=  1.029218231018771e-002  =H(37)
H( 5)= -1.031352505876851e-002  =H(36)
H( 6)= -6.867711169248183e-003  =H(35)
H( 7)=  1.800130147493761e-002  =H(34)
H( 8)= -7.781179501205805e-003  =H(33)
H( 9)= -1.593497517309733e-002  =H(32)
H(10)=  2.419735516010280e-002  =H(31)
H(11)= -1.371979800557242e-003  =H(30)
H(12)= -3.035606772323000e-002  =H(29)
H(13)=  2.979843200753757e-002  =H(28)
H(14)=  1.356420391493984e-002  =H(27)
H(15)= -5.512175099812947e-002  =H(26)
H(16)=  3.405661353296577e-002  =H(25)
H(17)=  5.370607795481440e-002  =H(24)
H(18)= -1.227353345368210e-001  =H(23)
H(19)=  3.634698998054064e-002  =H(22)
H(20)=  5.319161272267646e-001  =H(21)

```

Figure 5.19: Gabarit du filtre passe-bas FIR

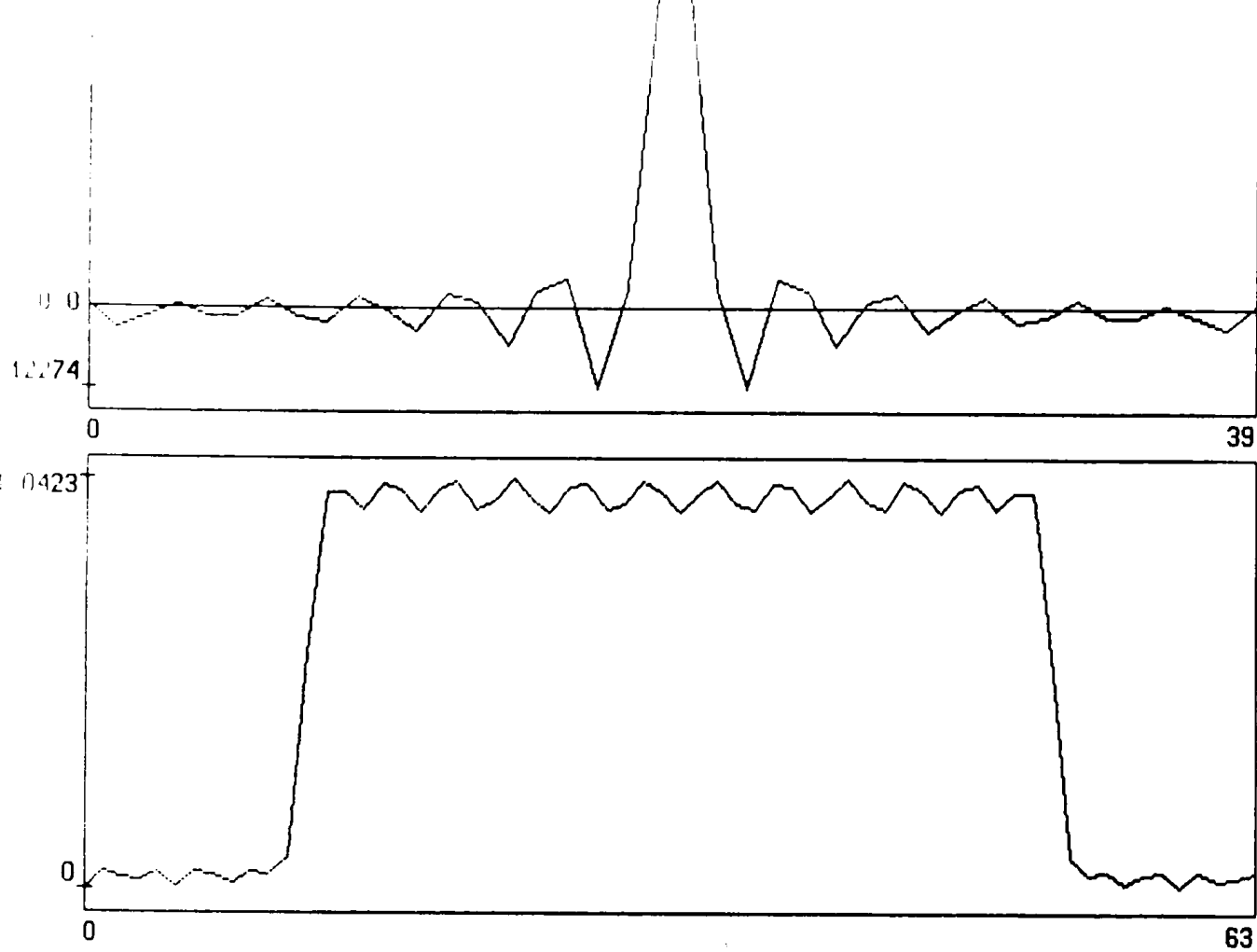


Figure 5.20 : réponses a)impulsionnelle et b)fréquentielle du filtre passe-bas FIR

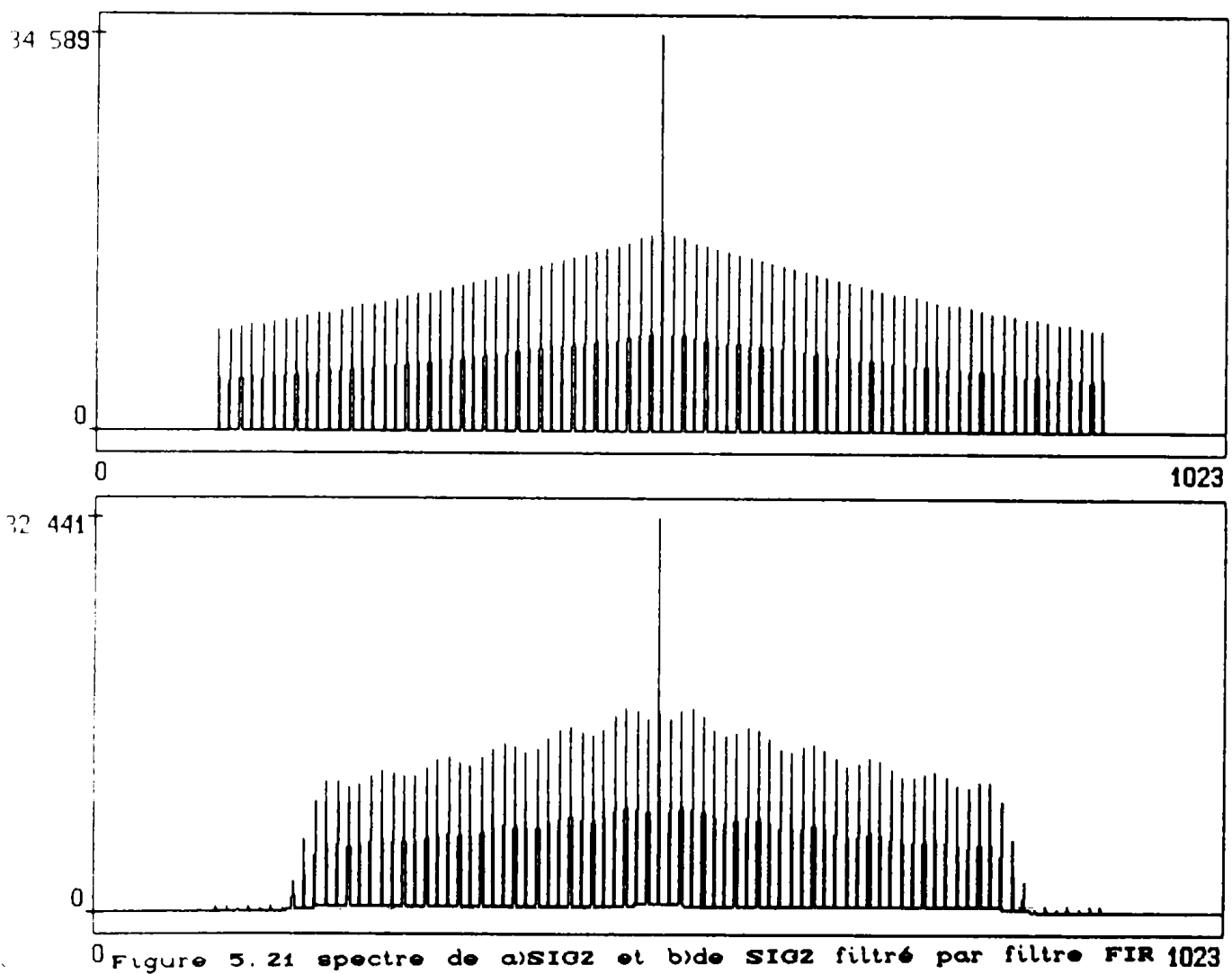


Figure 5.21 spectre de a)SIG2 et b)de SIG2 filtré par filtre FIR 1023

FINITE IMPULSE RESPONSE FOR A
LINEAR PHASE FILTER

BAND	LOWER BAND EDGE	UPPER BAND EDGE	DESIRED VALUE
1	0.000000e+000	3.300000e-001	1.000000e+000
2	3.600000e-001	5.000000e-001	0.000000e+000

FILTER LENGTH = 40

IMPULSE RESPONSE

H(1)= -2.246093750000000e-002 =H(40)
H(2)= 1.119995117187500e-002 =H(39)
H(3)= 6.469726562500000e-003 =H(38)
H(4)= -8.728027343750000e-003 =H(37)
H(5)= 1.480102539062500e-002 =H(36)
H(6)= 2.929687500000000e-003 =H(35)
H(7)= -1.275634765625000e-002 =H(34)
H(8)= 2.096557617187500e-002 =H(33)
H(9)= -2.685546875000000e-003 =H(32)
H(10)= -1.626586914062500e-002 =H(31)
H(11)= 3.067016601562500e-002 =H(30)
H(12)= -1.245117187500000e-002 =H(29)
H(13)= -1.913452148437500e-002 =H(28)
H(14)= 4.736328125000000e-002 =H(27)
H(15)= -3.213500976562500e-002 =H(26)
H(16)= -2.117919921875000e-002 =H(25)
H(17)= 8.731079101562500e-002 =H(24)
H(18)= -9.497070312500000e-002 =H(23)
H(19)= -2.221679687500000e-002 =H(22)
H(20)= 5.634155273437500e-001 =H(21)

Figure 5.22: Gabarit du filtre passe-bas FIR pour le signal EEG

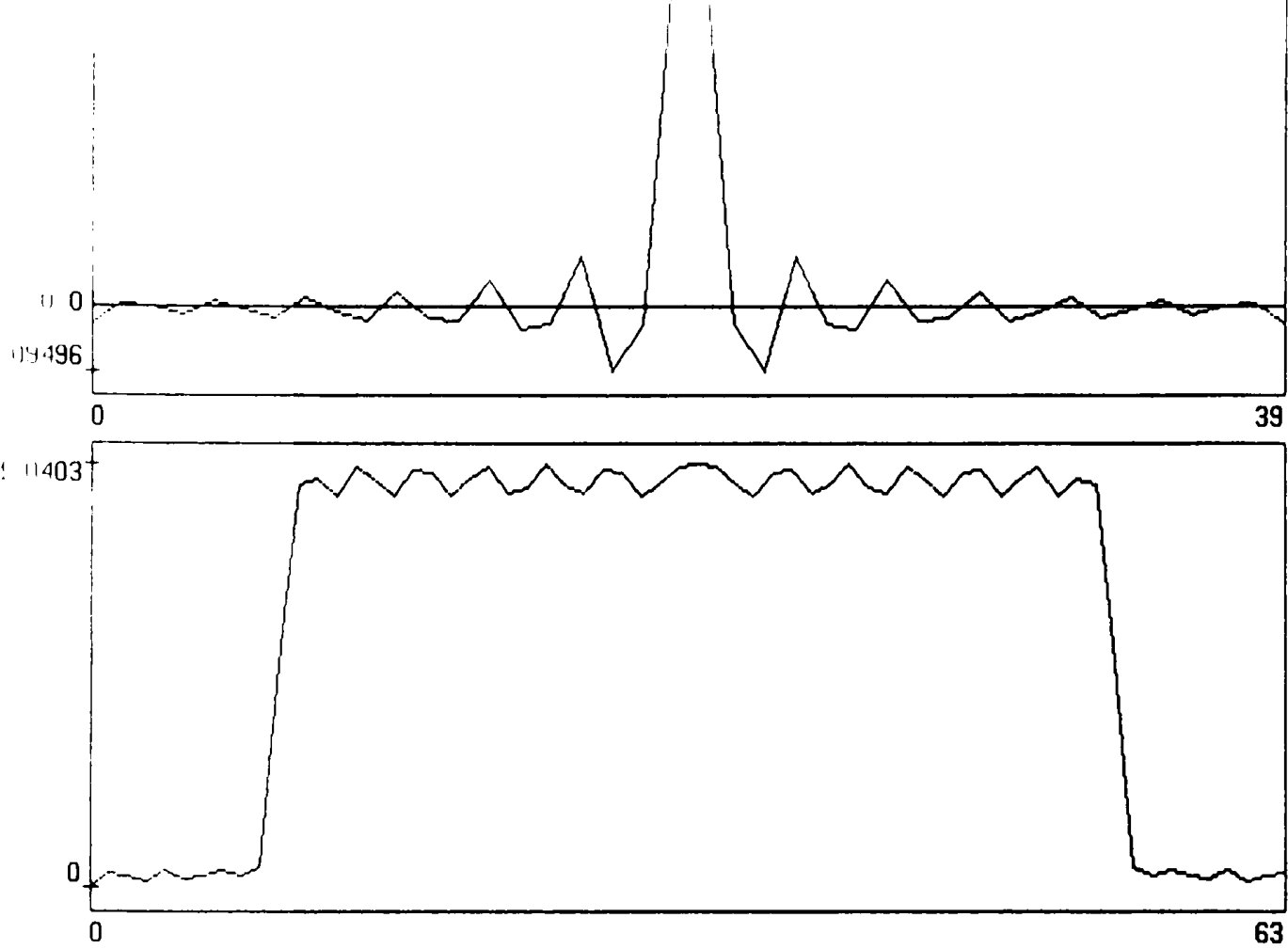


Figure 5.23 : a) réponse impulsionnelle du filtre passe-haut FIR pour le signal EEG b) réponse fréquentielle du filtre

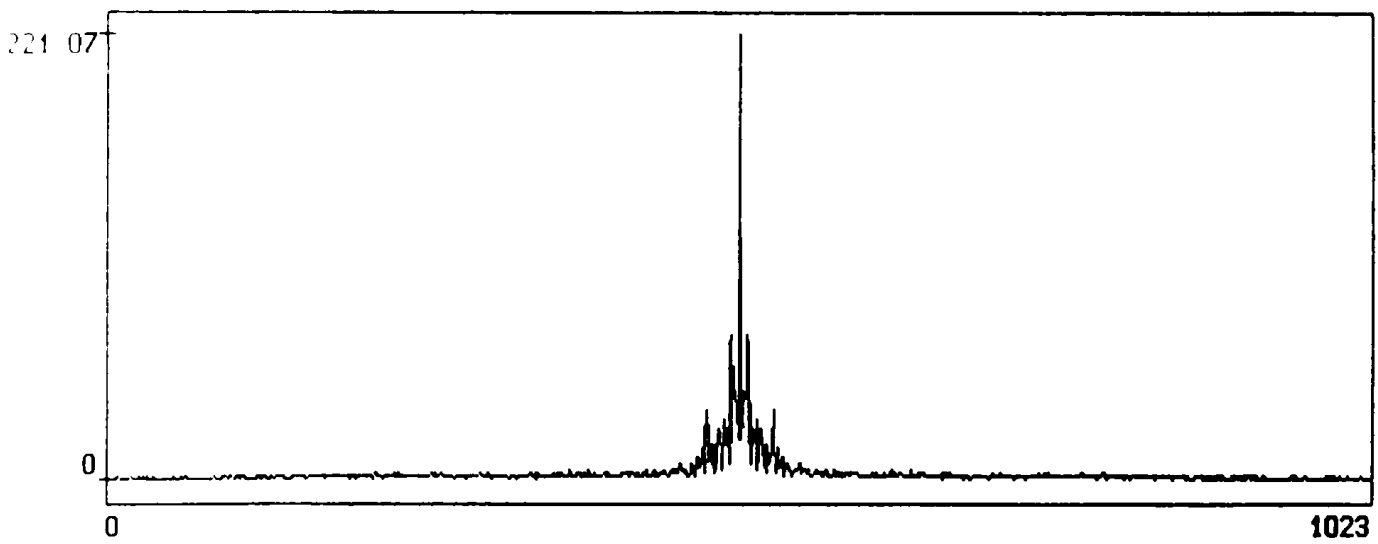
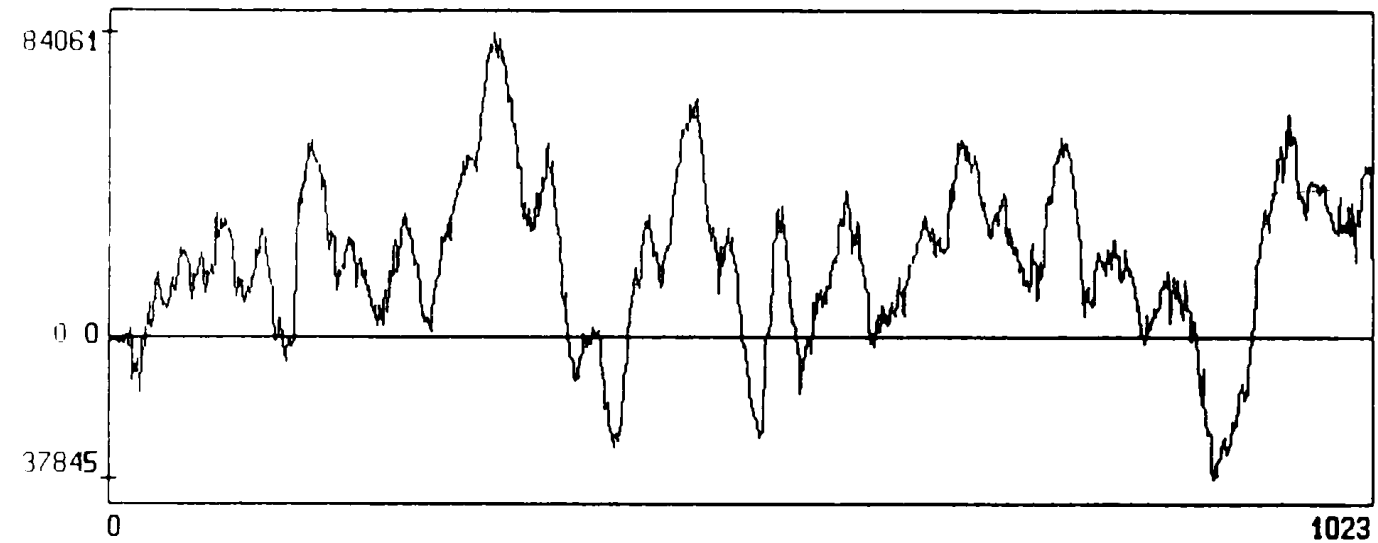
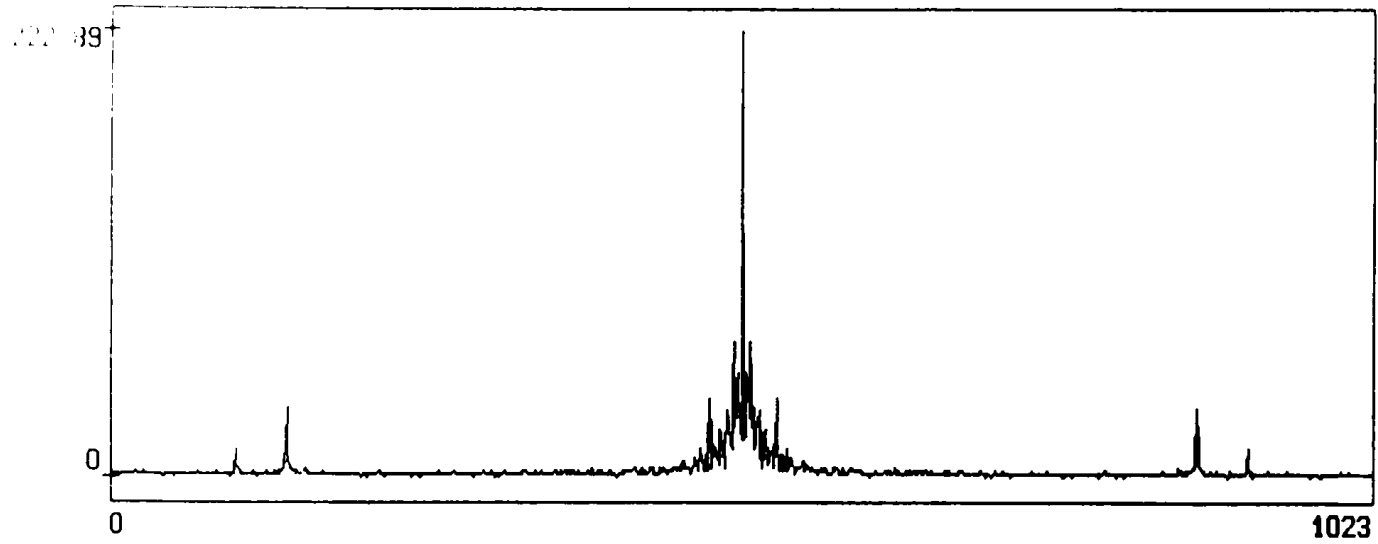
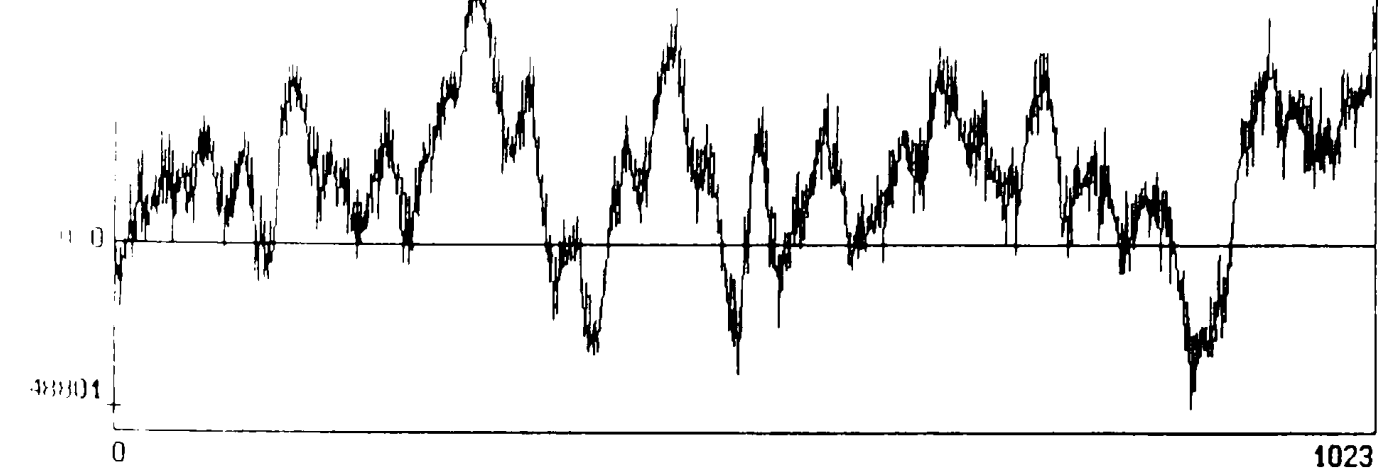


Figure 5.24 : a) signal EEG bruité b) spectre du signal EEG bruité
 c) signal EEG filtré par le filtre FIR d) spectre du EEG filtré

5.5 Estimation de la fonction d'autocorrélation du bruit de quantification :

Dans les chapitres précédents, on a supposé que l'erreur due à la quantification des signaux est un bruit blanc. On s'est proposé de vérifier ce résultat. On a pris 20 réalisations différentes de l'erreur de quantification d'un signal $x(n)$ selon la figure (5.25). On a ensuite calculé l'autocorrélation de chaque réalisation.

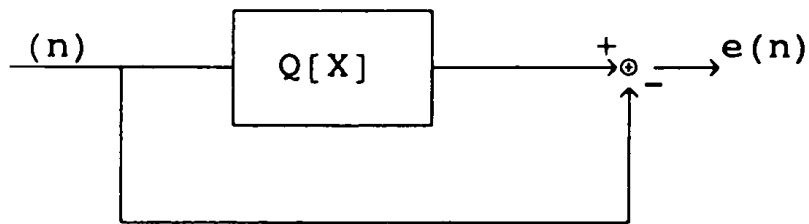


Figure 5.25 : Détermination de l'erreur de quantification

L'autocorrélation est calculée avec le programme réalisant le filtrage FIR. Les échantillons du signal $e(n)$ sont les coefficients du filtre et les échantillons du signal $e(-n)$ constituent le signal d'entrée. L'autocorrélation du processus est obtenu en faisant la moyenne des fonctions d'autocorrélation de chaque réalisation. Les figures (5.26) et (5.27) illustrent les deux premières réalisations avec leur autocorrélation respective. La figure (5.28) montre la moyenne des fonctions d'autocorrélation de 10 et de 20 réalisations.

De ces graphes, on voit qu'on s'approche de plus en plus vers l'autocorrélation du bruit blanc en augmentant le nombre de réalisations. Et on peut affirmer que l'erreur due à la quantification des signaux est un bruit blanc.

La fonction d'autocorrélation qui apparaît sur les graphes est normalisée par rapport à la valeur maximale.

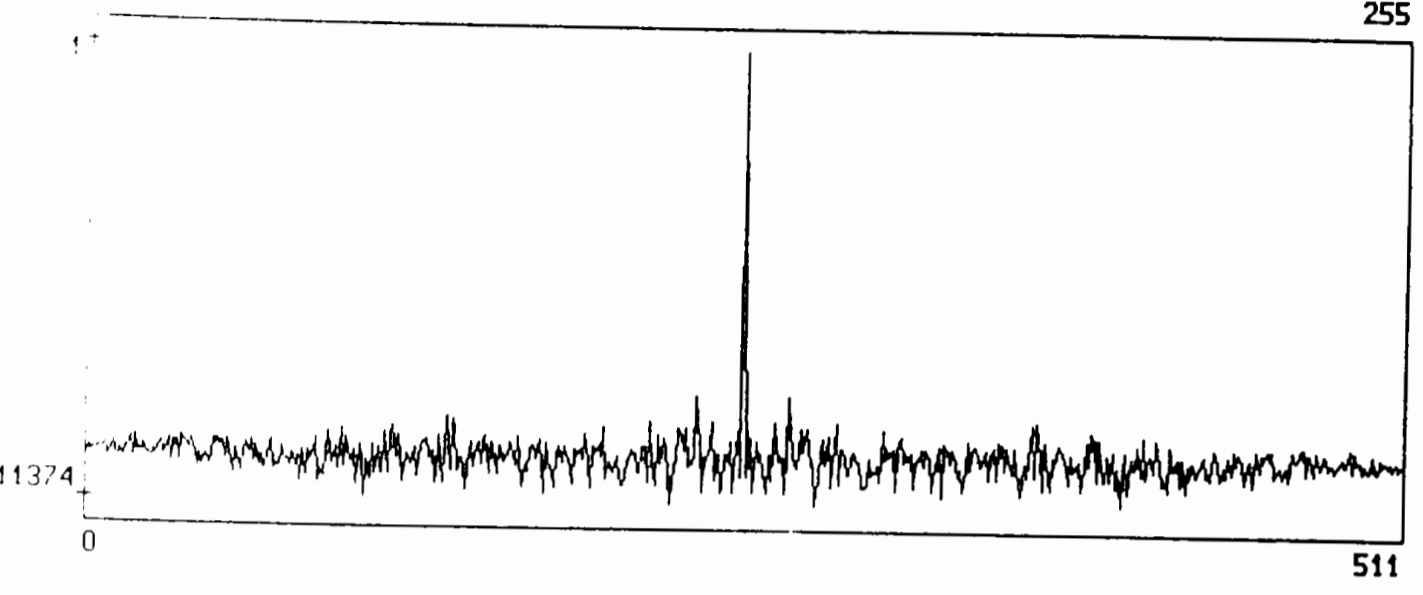
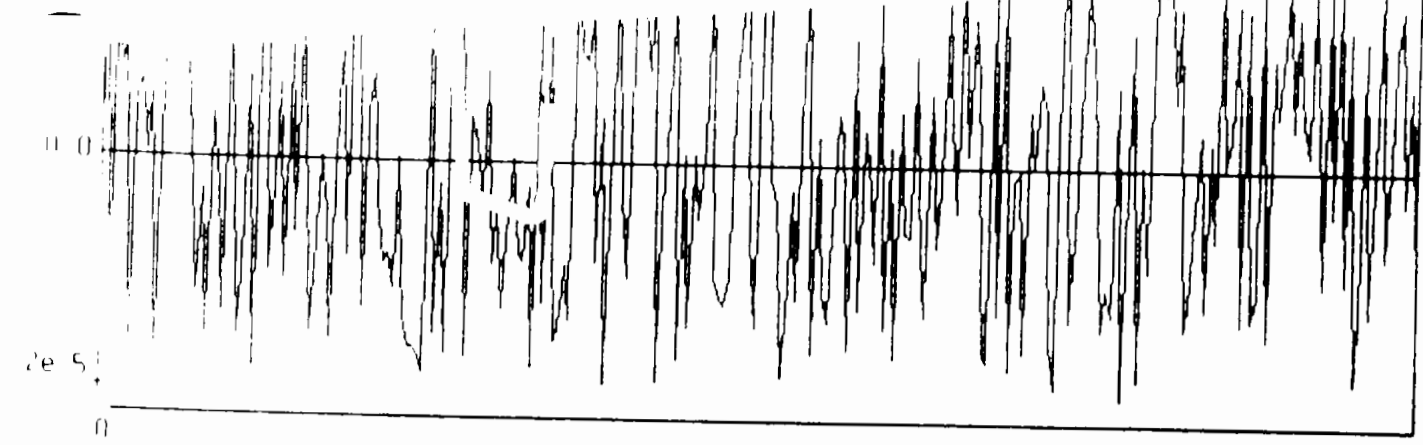


Figure 5.26: a) 1^{ère} réalisation du bruit et b) son autocorrélation

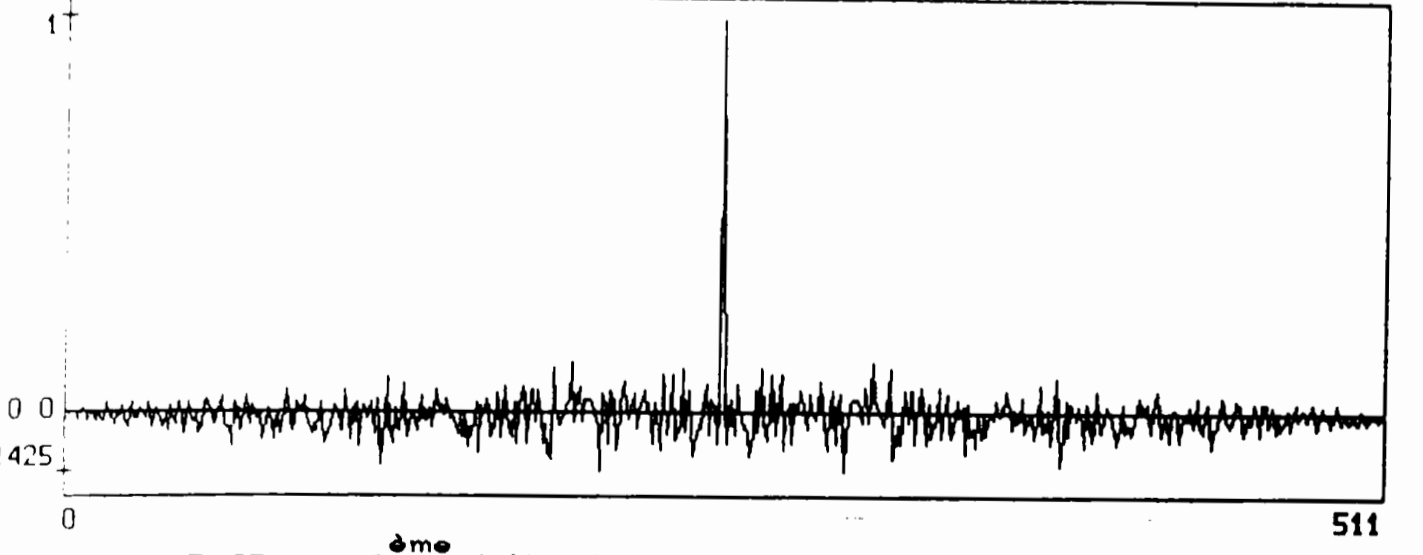
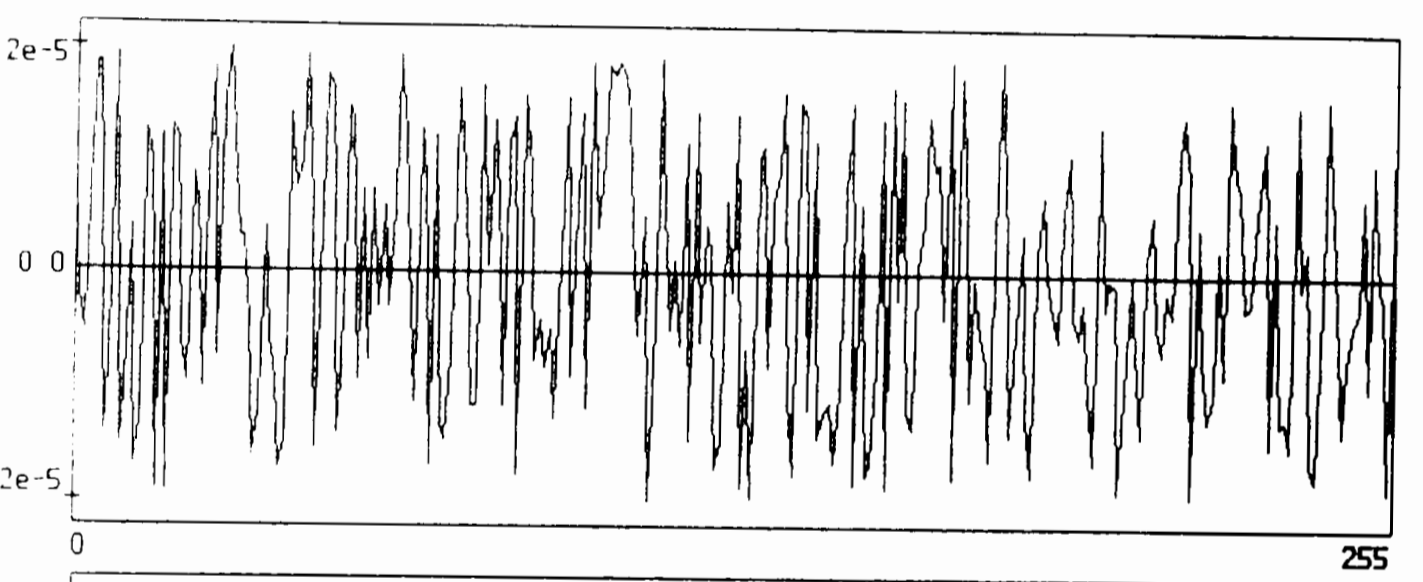


Figure 5.27: a) 2^{ème} réalisation du bruit et b) son autocorrélation

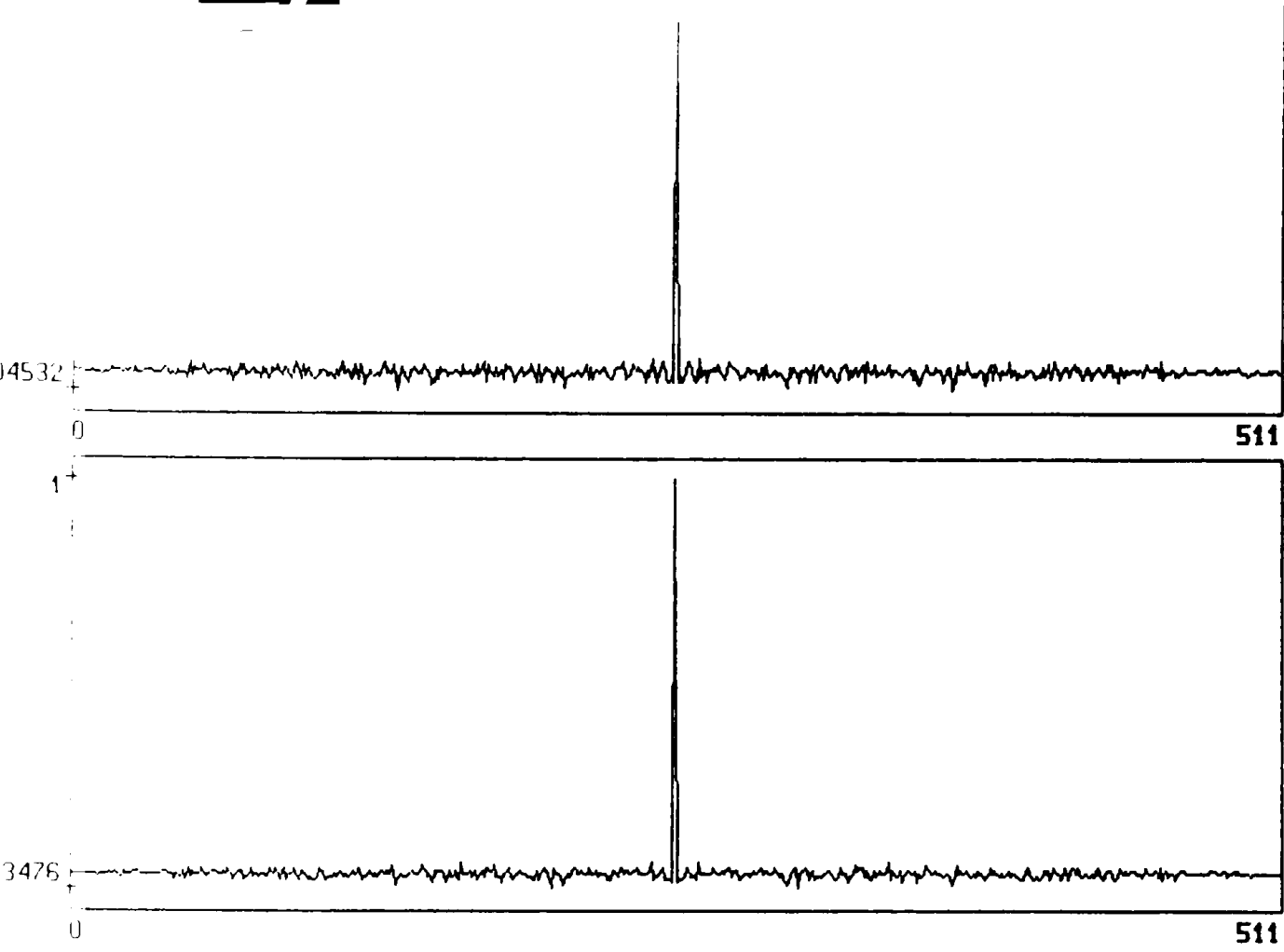


Figure 5.28: Autocorrélation de 10 et de 20 réalisations du bruit

CONCLUSION

L'objectif qu'on s'était fixé au début du travail a été atteint. Sur le plan théorique, on a étudié les algorithmes de base du traitement numérique du signal ainsi que les problèmes qui surviennent quand on les implante à l'aide d'une arithmétique à précision finie. On a ensuite mis au point ces algorithmes sur le système de développement d'un processeur de signal en les optimisant du point de vue compromis entre le temps de calcul et précision numérique.

Le travail est intéressant du point de vue didactique et de recherche. En effet, il est dorénavant possible aux étudiants de concevoir des applications avec une idée précise sur le temps de calcul et le bruit de quantification engendré. Du côté recherche, l'utilisateur dispose d'un noyau autour duquel il peut concevoir ses applications.

Les perspectives peuvent prendre les axes suivants:

- utiliser un processeur avec des signaux réels et concevoir des applications typiques tels que le traitement de la parole et le traitement de l'image.

- implanter et étudier les mêmes algorithmes sur le même processeur en virgule fixe multiprécision et en virgule flottante. On devra comparer les performances résultantes avec les résultats qu'on a obtenu dans le présent travail.

BIBLIOGRAPHIE

- [1] J.Allen, " Computer architecture for signal processing, " Proc.IEEE, vol.63, Apr.1975, pp 624-633.
- [2] H.J.Kolb, " Effective programming and realization of real time signal processors, " Signal processing : Theorie and Applications, 1980, pp 359-362.
- [3] Y.S.Wu, " Architectural considerations of a signal processor under microprogram control, " Spring Joint Comput.Conf., AFIPS conf.proc, Vol.40, May 16-18, 1972, pp. 675-683.
- [4] J.V.Harshman, " Architecture of a programmable Digital Signal Processor, " Nat-Telecommun.Conf.Rec, Dec.2-4, 1974, pp 496-500.
- [5] J.R.BODDIE, G.T.DARYANANI, U.ELDUMIATI, R.N.GADNEZ, J.S.THOMPSON, et S.M.Walters, " Architecture and Performance, " BSTJ, vol.60, Sep.1981, pp 1499-1462.
- [6] D.J DeFatta, J.G.Lucas et W.S.Hodgkiss, " Digital Signal Processing : A system Design Approach," John Willey and Sons, New-york, 1988.
- [7] A.V.Oppenheim et R.W.Schafer, " Digital Signal Processing, " Prentice-Hall. Inc, Englewood Cliffs, New Jersey, 1975.
- [8] M.Kunt, " Traitement numérique du signal, " Dunod, Editions Georgi, 1981.
- [9] N.B.Jones, " Digital Signal Processing, " Peter Peregrinus Ltd, UK, 1982.

- [10] A.PELED et B.Liu, " Digital Signal Processing : Theory, Design and implementation, " John Wiley and Sons, New York, 1976.
- [11] R.Boite et H.Leich , " Les filtres numériques, " Masson, Paris, 1982.
- [12] J.W.Cooley et J.W.Tuckey, " An algorithm for the Machine Computation of complex Fourier Series, " Math.Computation, vol 19, Apr.1965, pp.297-301.
- [13] F.J.Harris, " On the Use of Windows for Harmonic Analysis with the Discrete Fourier Transform," Proc.IEEE, vol.66, Jan 1978, pp.51-83.
- [14] A.Antoniou, " Digital Filters : Analysis and Design," McGraw-Hill, New York, 1979.
- [15] M.Bellanger, " Traitement numérique du signal : th+orie et pratique," Masson, Paris, 1987.
- [16] K.Steiglitz, " Computer-Aided Design of Recursive Digital Filters," IEEE Trans.Audio Electroacoust., vol. AU-18, June 1970.
- [17] L.R.Rabiner et K.Steiglitz, "The Design of Wide-Band Recursive and Non-Recursive Digital Differentiators," IEEE Trans Audio Electroacoust, vol.18, June 1970, pp.204-209.
- [18] A.G.Desczky," Synthesis of Recursive Digital Filters Using the minimum P Error Criterion," IEEE Trans.Electroacoust vol.AU-20, Oct 1972, pp 257-263.
- [19] P.Fondaneche et P.Gilbertas, "Filtres numériques : principes et réalisations," Masson, Paris, 1981.
- [20] W.D.Stanley, " Digital Signal Processing, " Prentice-Hall Compagny, Reston, Virginia, 1975.

- [21] R.E.Bogner et A.G.Constantinides, " Introduction to Digital Filters," Jhon-Wiley and Sons, Chichester, 1975.
- [22] L.R.Rabiner et B.Gold, " Theory and Applications of Digital Signal Processing, " Prentice-Hall, Englewood Cliffs, N.J, 1975.
- [23] T.W.Parks et J.H.McClellan, " Chebyshev Approximation for Non-Recursive Digital Filters with Linear Phase," IEEE trans on circuit theory, vol.CT-19, MAR 1972, pp.189-199.
- [24] J.H.McClellan, T.W.Parks et L.W.Rabiner," A Computer Program for Designing Optimum FIR Linear-Phase Digital Filters, "IEEE trans. on Audio Electroacoust, vol AU-21, Dec 1973, pp.506-526.
- [25] L.W.Rabiner, J.H.McClellan et T.W.Parks, " FIR Digital Filter Design Techniques Using Weighted Chebyshev Approximation," Proc IEEE, vol 63, Apr 1975, pp.595-610.
- [26] J.H.McClellan et T.W.Parks, " A United Approach to the Design of Optimum FIR Linear-Phase Digital Filters, " IEEE trans. on Circuit Theory, vol CT-20, Nov 1973, pp 697-701.
- [27] A.Papoulis, "Probability, Random Variables and Stochastic Process , " McGraw-Holl Book Compagny, New York, 1965.
- [28] Spataru, " Théorie de la transmission de l'informantion," Mason et Cie, 1970.
- [29] F.J.Taylor," Digital Filter Design Handbook," Marcel Dekker inc., 1983.
- [30] C.L.Philips et H.T.Nagle, " Digital Control System Analysis and Design," Prentice-Hall, Englewood Cliffs, N.J, 1984.
- [31] A.V.Oppenheim et C.J.Weinstein, " Effects of Register Length in Digital Filtering and the Fast Fourier Transform," Proc IEEE, vol 60, Aug 1972, pp 957-976.

- [32] L.B.Jackson, " On the Interaction of Roundoff Noise and Dynamic Range in digital Filters," BSTJ, vol 49, Feb 1970 pp.159-184.
- [33] R.Boite, " Les filtres numériques," Journ+e d'Electronique, Press polytechniques Romandes, Lausanne, 1981, pp 85-104.
- [34] L.B.Jackson, " Roundoff-Noise Analysis for Fixed-Point Digital Filters Realized in Cascade or Parallel Form," IEEE trans on Audio and Electroacoust., vol AU-18, June 1970, pp 107-122.
- [35] ADSP-2100 User's Manual, Analog Devices Inc, 1986.
- [36] M.Bouamar, " Fast DSP-Based System for topographic EEG Analysis, " Microprocessors and Microsystems, vol 15, Apr 1991, pp 160-166.
- [37] ADSP-2100 Cross-Software Manual, Analog Devices Inc, 1986.
- [38] ADSP-2100 Applications Handbook, vol 1 Analog Devices Inc., 1987.
- [39] Monarch User's Manual, Antenna-Group, Inc.