

10/91

وزارة الجامعات
Ministère aux Universités

ECOLE NATIONALE POLYTECHNIQUE

DEPARTEMENT ELECTRONIQUE

المدرسة الوطنية المتعددة التقنيات
BIBLIOTHEQUE — المكتبة
Ecole Nationale Polytechnique

PROJET DE FIN D'ETUDES

SUJET

SIMULATION DES FILTRES NUMERIQUES
R. I. I. EN VIRGULE FIXE

Proposé par : B. DERRAS

Etudié par : N. SADI AHMED

Dirigé par : B. DERRAS

PROMOTION

Juin 1991

Ecole Nationale Polytechnique

DEPARTEMENT D'ELECTRONIQUE

المدرسة الوطنية المتعددة التقنيات
BIBLIOTHEQUE — المكتبة
Ecole Nationale Polytechnique

PROJET DE FIN D'ETUDE

SIMULATION DES FILTRES NUMERIQUES RII EN VIRGULE FIXE

PROPOSEE ET DIRIGE PAR : M. B. DERRAS

ETUDIE PAR : SADI AHMED NAFISSA

PROMOTION JUIN 1991

Mes Remerciements . . .

*les plus chaleureux sont adressés à Mr B.DERRAS ,
mon promoteur, pour avoir fait preuve de générosité
et de compréhension à mon égard tout au long de
cette année . Je tiens aussi à lui exprimer ma
profonde gratitude pour le soutien et l'aide
permanents dont il m'a fait part afin de mener à bien
ce présent projet. De même, je saisis l'opportunité
de ces quelques lignes pour lui témoigner mon immense
respect.*

*Mes plus vifs remerciement s'adressent à M^{elle} Djamilia,
ingenieur au centre de calcul de l'Ecole pour son aide
précieuse .*

إهداء

إلى العبيبة أسمى ... إلى الغالي أسمى ...

إلى جزيرة و هناء و صلاح الدين ...

إلى كل من أحبهم ...

إلى كل قارئ شغوف بقصر الله سفاوة بما تعلمه هذه الصفحات

و أتمنى لمن يصيب سبغاه و سراره ...

نقيس

ملخص المشروع

الغاية الأساسية من هذه الدراسة هي التمكن من مضاهاة عمل المرشحات العددية ذات الإستجابة الدافعية اللامتناهية في الزمن ، بهدف إظهار مدى تأثير عملية الترشيح من جراء إستعمال سجلات محدودة الطول لحفظ المعطيات و إنجاز العمليات الحسابية ، و بخاصة في حالة التمثيل الرقمي بالفاصلة الكسرية الثابتة . كما تمكن هذه الدراسة من التحقق من صلاحية الفرضيات المستعملة في نظرية البنيات الفضلى عن طريق النتائج العملية ، بالإضافة إلى دراسة بنيتين ناتجتين عن تحليل المرشح إلى خلايا من المرتبة الثانية موصولة على التوازي ، و التي تشكل حلا وسطا بين البنيات الفضلى و المباشرة الإجمالية ، من حيث قيمة الخطأ الحسابي و بساطة الأجهزة الضرورية لتحقيقها . و قد يفسح هذا التصنيف مجالا أوسع لإختيار البنية المناسبة حسب المقاييس التقنية و الاقتصادية المحددة .

PROJECT ABSTRACT

This work is devoted to the simulation of IIR digital filters in fixed point arithmetic in order to analyse the output error behavior with the respect to the finite wordlength effects.

It has been shown that a minimum output error variance can be reached by the use of the digital state-space descriptions. But, unfortunately these structures present a high hardware complexity which leads to look for other structures presenting a compromise between noise gain and complexity.

Two types of compromise structures based on second order parallel section filter decomposition are proposed in this project in the aim to permit a larger range choice of the convenient structure for specific applications.

RESUME DU PROJET

La finalité essentielle de ce projet est de simuler les filtres numériques RII en virgule fixe afin d'analyser leur comportement vis à vis des effets de la limitation de la précision de toutes les valeurs intervenant dans l'opération du filtrage et de vérifier la validité des hypothèses faites dans la théorie des structures à gain minimal.

Deux types de structures basées sur la théorie de la décomposition du filtre en cellules du second ordre disposées en parallèle, sont étudiées et simulées. L'avantage de ces deux structures par rapport aux structures optimale et canonique est le bon compromis qu'elles offrent entre la qualité du filtrage et la complexité des circuits.

	PAGE
CHAPITRE I	
Introduction Générale.....	1
1-1 préliminaire.....	1
1-2 Organisation du projet.....	2
CHAPITRE II	
Généralités sur la synthèse des filtres numériques.....	4
2-1 Systèmes et signaux discrets.....	4
2-1-1 Définitions.....	4
2-1-2 Transformée de Fourier.....	5
2-1-3 Echantillonnage et recouvrement.....	6
2-1-4 Transformée en Z.....	6
2-1-5 Réponse Impulsionnelle.....	7
* Stabilité des systèmes discrets.....	7
2-1-6 Equation aux différences finies.....	8
2-1-7 Réponse fréquentielle.....	8
2-2 Les filtres numériques.....	9
A) Filtres à reponses impulsionnelles de durée finie (RIF).....	10
1- Méthode du fenêtrage temporel.....	10
2- Méthode d'échantillonnage en fréquence.....	11
3- Méthode du minimax.....	12
B) Filtres à reponses impulsionnelles de durée infinie (RII).....	12
1- Méthodes classiques.....	13
a- Méthode de l'invariance de la réponse impulsionnelle.....	13
b- Méthode de la transformation bilinéaire.....	14
2- Méthodes d'optimisation.....	15
a- Approximation de PADE.....	15
b- Méthode des moindres carrés.....	16
2-3 Transformation fréquentielle.....	17
2-4 Tableau comparatif entre les filtres RIF et RII.....	19
CHAPITRE III	
Effets de la longueur limitée des mots sur le filtrage numérique.....	20
3-1 Introduction.....	20
3-2 Conversion analogique numérique.....	21
3-3 La quantification.....	22
3-4 Types de représentation arithmétique binaire.....	23
3-4-1 Représentation en virgule fixe.....	24
3-4-2 Représentation en virgule flottante.....	25
3-5 Erreurs de calcul et de dépassement.....	25
3-5-1 Oscillations de dépassement.....	27
CHAPITRE IV	
Représentation des filtres numériques par les variables d'état.....	29
4-1 Introduction.....	29
4-2 Représentation d'état des filtres numériques.....	29
* Changement de coordonnées.....	31
4-3 Matrice K et la stabilité de LYAPUNOV.....	32
4-4 Structures éliminant les oscillations	

4-4	Structures éliminant les oscillations de dépassement.....	33
4-5	Normalisation des filtres numériques en virgule fixe.....	35
4-5-1	Règles de normalisation.....	35
4-5-1	Normalisation des paramètres d'état...	37
4-6	Analyse du bruit de calcul dans un filtre numérique RII.....	38
4-7	Remarque importante.....	42
CHAPITRE V	Structures optimales.....	43
5-1	Minimisation du gain de bruit de calcul dans le cas de registres de longueurs égales....	43
A)	Méthode de MULLIS-ROBERTS.....	43
B)	Méthode de HWANG.....	45
5-2	Algorithmes de calcul des matrices K et W...	48
5-3	Propriétés des structures à gain de bruit minimal.....	49
5-3-1	Invariance par rapport à une transformation fréquentielle.....	49
5-3-2	Expression complète du bruit de calcul.....	49
5-3-3	Gain minimal et performance des réalisations.....	50
5-4	Minimisation de la sensibilité aux coefficients.....	50
CHAPITRE VI	Structures décomposées optimales.....	53
6-1	Introduction.....	53
6-2	Structures optimales du second ordre.....	54
1-	Structures de BOMAR	55
2-	Structures de BARNES.....	56
3-	Structures de HWANG.....	58
6-3	Conclusion.....	59
CHAPITRE VII	Resultats et Interprétations.....	61
7-1	Introduction	61
7-2	Performances d'un filtre numérique avec une structure canonique.....	62
7-2-1	Simulation avec précision infinie.....	62
7-2-2	Simulation avec précision finie.....	63
7-2-3	Variance de l'erreur de sortie.....	63
7-2-4	Sensibilité de la réponse fréquentielle.....	64
7-3	Performances d'un filtre numérique avec une structure décomposée en parallèle.....	65
7-3-1	Simulation avec précision finie.....	68
a-	Variance de l'erreur de sortie.....	68
b-	Sensibilité de la réponse fréquentielle.....	69
7-4	Effets d'une transformation fréquentielle sur le gain de bruit de calcul.....	70
CHAPITRE VIII	Conclusion générale.....	71
	Bibliographie.....	73
	Annexe.....	

CHAPITRE I

INTRODUCTION GENERALE

CHAPITRE I

INTRODUCTION GENERALE

1-1 Préliminaire

Il y a quelques années, la numérisation des filtres analogiques était seulement un moyen pour simuler leurs comportements sur ordinateur, dans le but de tester leurs performances et de corriger leurs imperfections avant leurs réalisations analogiques définitives. Ce moyen est très efficace et économique, car il peut en quelques sortes assurer le bon fonctionnement des filtres qui constitue un facteur capital dans la plupart des applications.

De nos jours, le grand développement de la technologie des circuits intégrés numériques offre avec la complicité de l'arsenal mathématique du traitement numérique des signaux, la possibilité de réaliser et d'utiliser des filtres proprement numériques. Ces derniers se distinguent par leur grande fiabilité, par leurs vitesses de traitement et par leurs coûts réduits.

La simulation de ces filtres avant leur implantation permet de résoudre les problèmes que peut poser leur réalisation matérielle; en plus des problèmes classiques de l'échantillonnage, le facteur de qualité, la binarisation et la représentation arithmétique,... il y a le problème de la précision avec laquelle doit s'opérer le filtrage du fait de l'utilisation de registres de tailles réelles pour représenter les valeurs des paramètres des filtres et des signaux qui les parcourent ainsi que les valeurs résultant des diverses opérations mathématiques.

L'impact de la limitation de la longueur des mots sur la qualité du filtrage, et les effets néfastes qu'elle cause, méritent une sérieuse considération et une plus grande part d'attention.

Les nombreuses investigations qui ont été engagées dans ce sens ont permis de prévaloir l'efficacité de la représentation des filtres numériques dans l'espace d'état par l'opportunité de trouver des structures ayant des caractéristiques optimales vis-à-vis des effets de la longueur finie des mots.

Cependant, ces structures optimales présentent une certaine complexité dans leurs réalisations vu le nombre d'opérations à effectuer donc du coût de la réalisation et du temps de calcul.

Ce qui a amené à chercher des structures de compromis qui ont pour objectif de réduire, simultanément la variance du bruit de calcul par rapport à celle de la structure canonique et la complexité matérielle par rapport à la structure optimale.

Ce présent travail s'intéresse à la simulation des filtres numériques RII en virgule fixe pour mettre en relief les conséquences de l'utilisation de registres de longueur finie sur la qualité du filtrage et pour vérifier l'applicabilité de la théorie de minimisation du gain de bruit de calcul. Les structures optimales sont approchées par deux méthodes: **MULLIS-ROBERTS** et **HWANG**, en diminuant la probabilité de dépassement de la gamme des valeurs représentables par une normalisation convenable.

Deux structures de compromis sont proposées, elles sont basées sur la décomposition du filtre en cellules du second ordre mises en parallèles dans l'une, les cellules sont optimales séparément et dans l'autre, elles sont canoniques. Les performances de telles structures du point de vue gain de bruit de calcul, nombre d'opérations, la sensibilité de la fonction de transfert ainsi que l'invariance à une transformation fréquentielle sont illustrés par la simulation d'un filtre RII prototype.

1-2 Organisation du projet

Cette étude débute par une brève présentation des outils mathématiques nécessaires pour la caractérisation et la synthèse des filtres numériques ainsi que leurs propriétés (chapitre II).

Le chapitre III est consacré à l'analyse des effets dus à la limitation de la longueur des mots et à mettre en évidence le compromis qui existe entre le bruit de calcul et les erreurs de dépassement dans les filtres numériques RII, tout en présentant les différents types de quantification d'arithmétique et de caractéristique de dépassement.

Les chapitres V et IV sont consacrés aux méthodes et procédures d'optimisation des erreurs de calcul en sortie du filtre par des structures optimales et des structures de compromis qui sont présentées par des cellules du second ordre optimales ou canoniques connectés en parallèles.

Le chapitre VI contient les résultats de la simulation des différentes structures pour un filtre donné, ainsi que l'illustration de certains résultats théoriques.

Et enfin, vient la conclusion générale au chapitre VIII.

CHAPITRE II

**GENERALITES SUR LA SYNTHESE
DES FILTRES NUMERIQUES.**

CHAPITRE II

GENERALITES SUR LA SYNTHÈSE DES FILTRES NUMÉRIQUES

2.1 Systèmes et signaux discrets

2.1.1 Définitions

Un signal discret est une suite de valeurs numériques réelles ou complexes. Il est dit numérique si son amplitude est quantifiée. Les valeurs du signal sont, en général, régulièrement réparties dans le temps. L'intervalle de temps entre deux valeurs du signal discret est appelé "période d'échantillonnage".

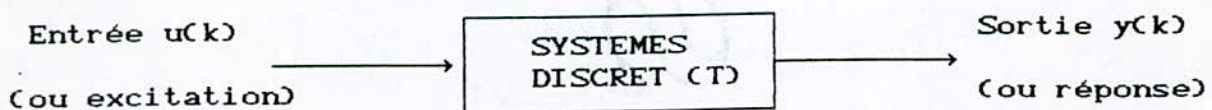
Comme dans le cas continu, les signaux discrets sont classés selon deux catégories:

- Les signaux déterministes dont l'évolution en fonction de la variable temps peut être prédite parfaitement par une représentation mathématique appropriée.

- Les signaux aléatoires, par contre, ont un comportement imprévisible; ils se caractérisent par leurs propriétés statistiques et fréquentielles.

Un signal périodique dont le comportement est imprévisible sur une période et qui est généré à partir d'une relation de récurrence parfaitement déterministe, est dit "pseudoaléatoire".

Un système discret agit sur un signal discret à son entrée pour produire un autre signal discret à sa sortie.



Fig(2-1) - Schématisation d'un système discret.

Un signal discret est causal s'il ne produit de réponse à un instant donné k , que pour les excitations antérieures; cela est dû au fait que, pour les systèmes physiques, l'effet ne peut précéder la cause.

Un système discret est linéaire, s'il réalise le principe de superposition; si T est l'opérateur caractérisant l'effet du système alors:

$$T[au(k) + bu'(k)] = a T[u(k)] + b T[u'(k)] \quad (2-1)$$

où a et b sont des constantes et u(k) et u'(k) sont deux entrées différentes à l'instant k (entier).

Un système discret est invariant dans le temps si :

$$T[u(k)] = y(k) \Rightarrow T[u(k-k_0)] = y(k-k_0) \quad (2-2)$$

où u(k) et y(k) sont respectivement l'entrée et la sortie du système T et k₀ un entier [1].

2.1.2 Transformée de FOURIER

La transformée de FOURIER (TF) d'un signal discret u(k) est donnée par :

$$U(f) = \sum_{k=-\infty}^{+\infty} u(k) e^{-j2\pi f k} \quad (2-3)$$

où f est la fréquence.

Une condition suffisante pour que U(f) existe est :

$$\sum_{k=-\infty}^{+\infty} |u(k)| < \infty \quad (2-4)$$

U(f) est une fonction complexe ou réelle de f et est périodique de période unité. En général, on considère U(f) sur l'intervalle principal [-1/2, 1/2]; La transformée inverse de U(f) est alors:

$$u(k) = \int_{-1/2}^{1/2} U(f) e^{j2\pi f k} df \quad (2-5)$$

Pour un signal u(k) réel, le spectre d'amplitude |U(f)| est paire et le spectre de phase arg[U(f)] = arctg[Im(U(f))/Re(U(f))] est impaire.

Si en plus u(k) est périodique de période N, on définit la transformée de FOURIER discrète (TFD):

$$U(n) = \sum_{k=0}^{N-1} u(k) e^{-2j\pi n k / N} \quad (2-6)$$

et donc

$$u(k) = \frac{1}{N} \sum_{n=0}^{N-1} U(n) e^{2j\pi nk/N} \quad (2-7)$$

Où $U(n)$ représente le domaine fréquentiel et $u(k)$ représente le domaine temporel [1],[2].

La TFD est très efficace dans le traitement numérique des signaux car c'est la seule transformation qui peut être programmées sur ordinateur; des algorithmes appelés transformée de FOURIER rapide (TFR) ont été développés pour accélérer et simplifier le calcul de la TFD [3].

2.1.3 Echantillonnage et recouvrement

L'échantillonnage consiste en l'observation d'un signal continu à certains instants qui sont généralement séparés par une durée constante T_e appelée période d'échantillonnage. Cette discrétisation dans le temps entraîne une répétition périodique du spectre du signal continu.

Pour que cette périodicité dans le domaine fréquentiel ne déforme pas le spectre répété (phénomène du recouvrement spectral) et permet la reconstitution du signal à partir de ses échantillons, il faut et il suffit que la fréquence d'échantillonnage (où de répétition) $F_e = 1/T_e$ soit égale ou supérieure au double de la fréquence maximale contenue dans le signal [4].

2.1.4 Transformée en Z

La transformation en Z est un outil plus puissant que la transformation de FOURIER qui ne représente qu'un cas particulier de celle-ci. Pour un signal $u(k)$, on définit sa transformée en Z bilatérale comme étant:

$$Z [u(k)] = U(Z) = \sum_{k=-\infty}^{+\infty} u(k) z^{-k} \quad (2-8)$$

Pour les systèmes et les signaux causals, on utilise la forme unilatérale:

$$U(z) = \sum_{k=0}^{+\infty} u(k) z^{-k} \quad (2-9)$$

Cette transformée existe pour un ensemble de points du plan complexe Z appelé "region de convergence", et dans laquelle:

$$\sum_k |u(k) z^{-k}| < \infty \quad (2-10)$$

La transformée en Z inverse est donnée par:

$$u(k) = \frac{1}{2\pi j} \oint_{\Gamma} U(z) z^{k-1} dz \quad (2-11)$$

où Γ est un contour fermé entourant l'origine dans le sens trigonométrique positif à l'intérieur de la région de convergence.

2.1.5 Réponse impulsionnelle

Un système linéaire est caractérisé complètement dans le domaine temporel par sa réponse impulsionnelle, qui est la réponse du système à une impulsion de Kronecker donnée par:

$$\delta(k) = \begin{cases} 1 & \text{si } k=0 \\ 0 & \text{sinon} \end{cases}$$

Pour une entrée $u(k)$, la sortie $y(k)$ d'un système discret s'exprime par:

$$y(k) = \sum_{l=-\infty}^{+\infty} u(l) h(k-l) \quad (2-12)$$

$$= (u * h)(k) = (h * u)(k)$$

où $h(k)$ est la réponse impulsionnelle du système et $*$ est l'opérateur de la convolution.

$$\text{Pour un système causal } h(k)=0 \quad \text{pour } k < 0 \quad (2-13)$$

* Stabilité des systèmes discrets

Un système linéaire discret est stable au sens de BIBO, si pour une entrée bornée, le système délivre une séquence bornée à sa sortie, autrement dit la norme de $h(k)$ est finie [1]:

$$\|h\|_1 = \sum_k |h(k)| < \infty \quad (2-14)$$

La transformée en Z de l'équation (2-12) est:

$$Y(Z) = H(Z) U(Z) \quad (2-15)$$

$$\text{où } H(Z) = \frac{Y(Z)}{U(Z)} \quad (2-16)$$

est appelée fonction du transfert du système.

On rappelle que l'expression (2-15) est obtenue en utilisant la propriété de convolution de la transformée en Z [1],[2].

2.1.6 Equation aux différences finies

L'excitation $u(k)$ et la réponse $y(k)$ d'un système numérique linéaire invariant dans le temps (LIT), satisfont une equation linéaire aux différences finies à coefficients constants de la forme:

$$\sum_{n=0}^N a_n y(k-n) = \sum_{m=0}^M b_m u(k-m) \quad (2-17)$$

ou encore

$$y(k) = \sum_{m=0}^M \frac{b_m}{a_0} u(k-m) - \sum_{n=1}^N \frac{a_n}{a_0} y(k-n) \quad (2-18)$$

où M , N , n et m sont des entiers. N est appelé l'ordre du système.

2.1.7 Réponse fréquentielle

En passant à la transformée en Z, on peut écrire à partir de l'equation (2-17) :

$$\sum_{n=0}^N a_n Y(z) z^{-n} = \sum_{m=0}^M b_m U(z) z^{-m} \quad (2-19)$$

La fonction de transfert du système s'exprime alors par:

$$H(z) = \frac{Y(z)}{U(z)} = \frac{\sum_{m=0}^M b_m z^{-m}}{\sum_{n=0}^N a_n z^{-n}} \quad (2-20)$$

La réponse fréquentielle du système est un cas particulier de la fonction de transfert $H(Z)$ pour laquelle $|z|=1$ (ou $z=e^{2j\pi f}$), elle est donnée par:

$$\begin{aligned} H(z) &= \frac{\sum_{m=0}^M b_m e^{-2j\pi f m}}{\sum_{n=0}^N a_n e^{-2j\pi f n}} \\ &= G \frac{\prod_{i=1}^M (e^{2j\pi f} - z_i)}{\prod_{i=1}^N (e^{2j\pi f} - p_i)} \end{aligned} \quad (2-21)$$

où G est une constante.

$H(f)$ caractérise complètement le système dans le domaine fréquentiel:

$$y(k) = h(k) * u(k) \xrightarrow{\text{TF}} Y(f) = H(f) U(f) \quad (2-22)$$

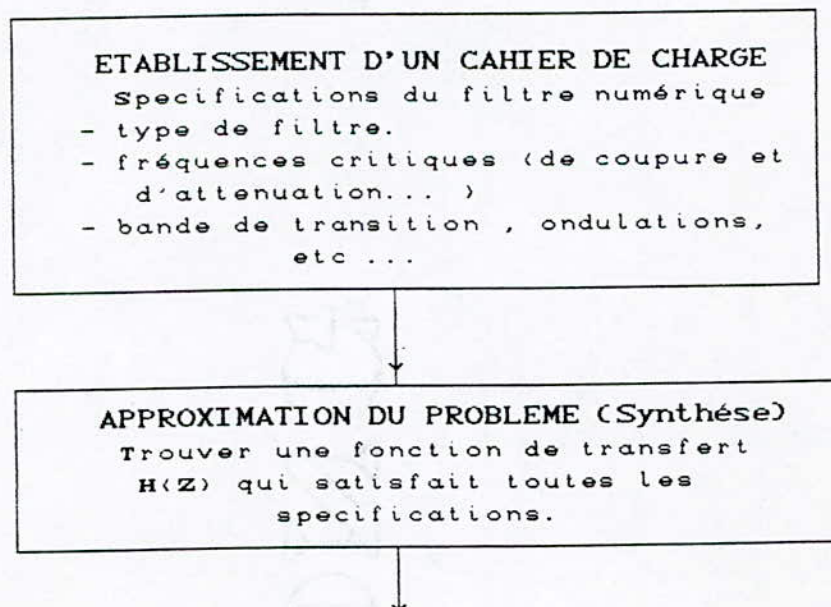
Les racines z_i du polynôme du numérateur de $H(f)$ sont appelées les zéros du système et les racines p_i du dénominateur de $H(f)$ sont appelées les pôles du système.

2.2 Les filtres numériques

Les filtres numériques sont des systèmes physiques donc causals qui ont pour fonction de modifier la distribution des composantes fréquentielles d'un signal selon des spécifications bien déterminées, en utilisant des opérations arithmétiques de précision limitée. On s'intéressera surtout aux systèmes LIT.

Ces filtres ont été conçus pour simuler les filtres analogiques sur ordinateur. Ce qui a permis d'étudier leurs performances et d'optimiser leurs paramètres avant leurs éventuelles réalisations. Mais avec le grand essor qu'a connu la technologie des circuits intégrés, l'intérêt des filtres numériques s'est accru et des méthodes de synthèse propres à ce type de filtres ont été largement développées.

La conception des filtres numériques se fait généralement suivant le schéma ci-dessous:



MODELISATION DU PROBLEME
 Simulation et mise au point
 Obtenir une réalisation qui définit
 la structure interne du filtre $H(Z)$.
 Cette réalisation doit être optimale vis-à-vis de la performance et des moyens actuels de calcul (ou traitement)

CONCRETISATION DU PROBLEME
 Détermination du matériel nécessaire pour réaliser le filtre avec une complexité réduite, une performance satisfaisante et un débit de données optimal.

On distingue deux grandes classes de filtres numériques :

A- Filtres à réponse impulsionnelle de durée finie (RIF)

Où $h(k)$ est non nulle sur un intervalle de temps N fini. Ils sont réalisés de manière non recursive suivant l'expression suivante:

$$y(k) = \sum_{n=0}^{N-1} h(n) u(k-n) \quad (2-17)$$

Leur schéma fonctionnel est de la forme [2]:

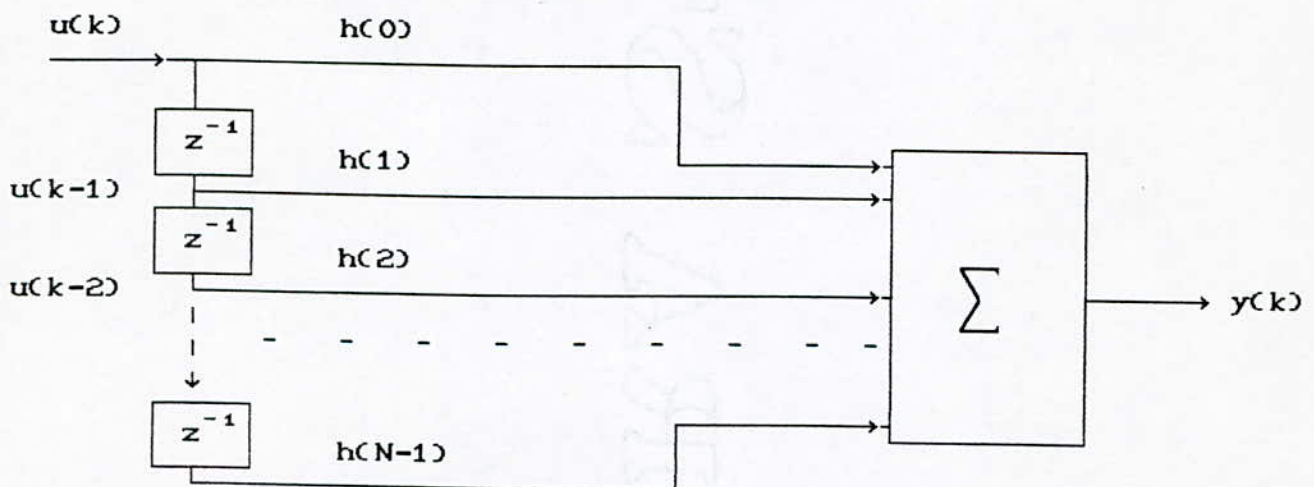


Fig (2-2) - Schéma fonctionnel d'un filtre RIF

On peut synthétiser ce type de filtres par plusieurs méthodes, parmi lesquelles :

1- Methodes du fenêtrage temporel

Soit un filtre de réponse impulsionnelle

$$h(k) = F^{-1}[H(\omega)]$$

non causale et infinie. On choisit une fenêtre $w(k)$ qui limite la durée de $h(k)$ dans $[k_1, k_2]$ et construit:

$$\hat{h}(k) = \begin{cases} h(k-k_2) w(k-k_2) & 0 \leq k \leq k_2 - k_1 = N-1 \\ 0 & \text{ailleurs} \end{cases} \quad (2-24)$$

qui est ainsi causale; on détermine alors le filtre RIF

$$\begin{aligned} \hat{H}(\omega) &= \sum_{k=0}^{N-1} \hat{h}(k) e^{-j\omega k} \\ &= e^{-j\omega(N-1)/2} \sum_{m=k_1}^{k_2} h(m) w(m) e^{-jm\omega} \end{aligned} \quad (2-25)$$

La fenêtre est choisie de telle sorte que:

- la linéarité de la phase soit conservée
- l'erreur

$$\epsilon^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |H(\omega) - \hat{H}(\omega)|^2 d\omega \quad (2-26)$$

soit minimale

- Le rapport (Amplitude du lobe principale)/(Amplitude du 1^{er} lobe secondaire) et la largeur du lobe principal (résolution de la fenêtre) soient acceptables pour conserver le maximum d'énergie de la réponse fréquentielle $H(\omega)$ à approximer [1].

2- Méthode d'échantillonnage en fréquence

Soit $H(\omega)$ la réponse fréquentielle spécifiée qu'on échantillonne en N points du cercle unité, puis on détermine la réponse impulsionnelle par la TFD inverse:

$$H(n) = \sum_{k=0}^{N-1} h(k) e^{-j2\pi nk/N} \quad (2-27)$$

$$h(k) = \frac{1}{N} \sum_{n=0}^{N-1} H(n) e^{j2\pi nk/N}$$

Ayant les N échantillons de $h(k)$, on calcule sa transformée en Z

$$\begin{aligned} H(z) &= \sum_{k=0}^{N-1} h(k) z^{-k} \\ &= \frac{1}{N} \sum_{n=0}^{N-1} H(n) \frac{1 - z^{-N}}{1 - z^{-1} e^{j2\pi n/N}} \end{aligned} \quad (2-28)$$

La réponse fréquentielle approximée $\hat{H}(\omega)$ est donc

$$\hat{H}(\omega) = H(z) \Big|_{z=e^{j\omega}} = \frac{1}{N} \sum_{n=0}^{N-1} H(n) I_N(n, \omega) \quad (2-29)$$

$$\text{où} \quad I_N(n, \omega) = \frac{1 - e^{-j\omega N}}{1 - e^{-j\omega} e^{j2\pi n/N}}$$

est appelée fonction d'interpolation [5].

Le nombre N des échantillons et les valeurs de $H(\omega)$ dans les bandes de transition seront choisis de façon à respecter les contraintes sur le niveau des ondulations et sur la largeur de la bande de transition [1].

3- Méthode du minimax

Cette méthode consiste à trouver une approximation $\hat{H}(e^{j\omega})$ de $H(\omega)$ de façon à minimiser la norme [1]:

$$\|\hat{H} - H\|_{\infty} = \max_{\omega} |\hat{H}(e^{j\omega}) - H(e^{j\omega})| \quad (2-30)$$

La méthode la plus classique est de trouver $H(e^{j\omega})$ sous forme d'un polynôme de CHEBYSHEV de variable:

$$x = \cos \omega \quad (-1 \leq x \leq 1) \quad (2-31)$$

et de degré au plus égal au nombre d'échantillons de $h(k)$ tel que

$$\text{si} \quad H(\omega) = p(x = \cos \omega), \quad H(\omega) = d(x) \quad \text{avec}$$

$$\varepsilon = \|p - d\|_{\infty} = \max_{-1 \leq x \leq 1} |p(x) - d(x)|$$

alors

$$\begin{cases} p(x_i) - d(x_i) = (-1)^i \varepsilon \\ (x_i^2 - 1) (p'(x_i) - d'(x_i)) = 0 \end{cases} \quad (2-32)$$

$$\text{pour} \quad -1 \leq x_i < x_{i+1} \leq 1 \quad i=0, 1, \dots$$

Cette méthode garantit des ondulations d'amplitudes égales (théorème de CHEBYSHEV) dans la bande passante et dans la bande d'affaiblissement selon l'algorithme utilisé [1],[6].

B- Filtrés à réponses impulsionnelles de durée infinie (RII)

Leurs réalisations sont récursives, car la sortie de ce type de filtres dépend non seulement des entrées précédentes et présentes, mais aussi des sorties antérieures. Leurs fonctions de transfert sont de la forme:

$$H(z) = \sum_{k=0}^{+\infty} h(k) z^{-k} = \frac{\sum_{m=0}^M b_m z^{-m}}{\sum_{n=0}^N a_n z^{-n}} \quad (2-33)$$

avec $M \leq N$ (causalité).

Les constantes a_i et b_i ($i=0,1,\dots$) sont appelées coefficients du filtre.

Le schéma fonctionnel pour de tels filtres est le suivant:[2]

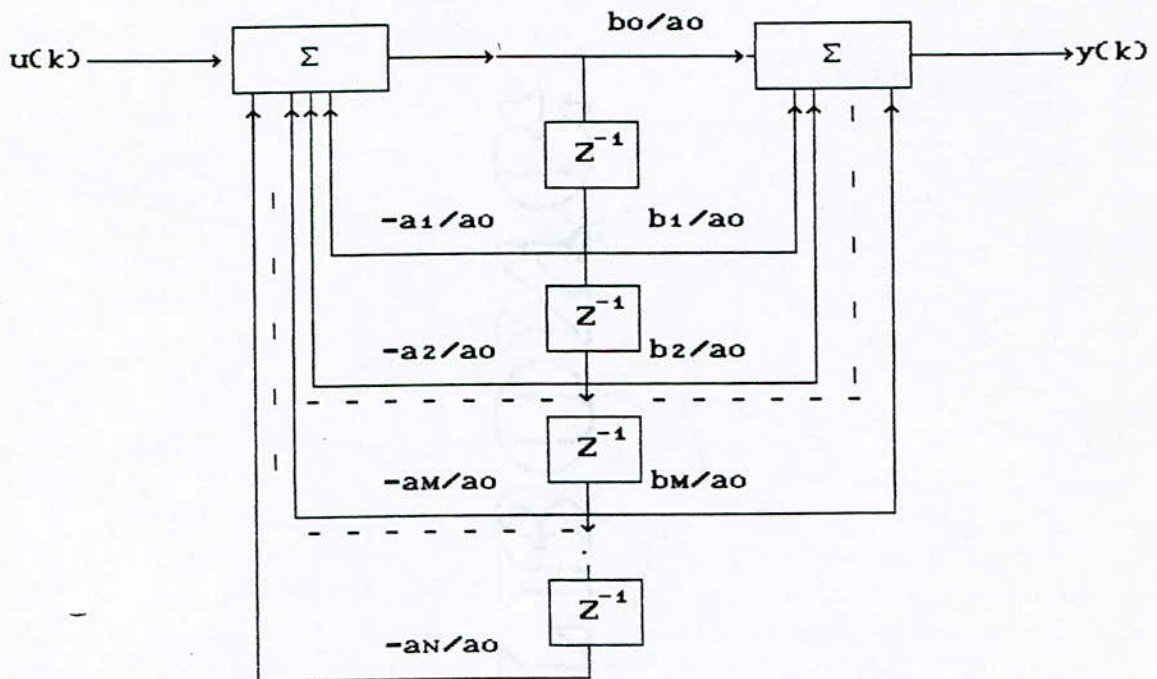


Fig (2-3)- Schéma fonctionnel d'un filtre RII

Les méthodes de synthèse de ce type de filtres se présentent en deux catégories:

- Les méthodes classiques qui établissent une correspondance entre les domaines analogique et numérique [2].
- Les méthodes algorithmiques d'optimisation assistées par ordinateur [6].

1- Méthodes classiques

a- Méthode de l'invariance de la réponse impulsionnelle

Le principe de cette méthode est de trouver une réponse impulsionnelle $h(k)$ d'un filtre numérique, qui soit la forme échantillonnée de la réponse approximée $\hat{h}(t)$ du filtre analogique désiré [2].

$$h(k) = T_e \hat{h}(k T_e)$$

où T_e est la période d'échantillonnage,

avec

$$\hat{H}(s) = L(\hat{h}(t)) = \frac{\alpha_0 + \alpha_1 s + \dots + \alpha_M s^M}{1 + \beta_1 s + \dots + \beta_N s^N} \quad (2-35)$$

où L est l'opérateur de la transformée de LAPLACE.

$$\text{donc } \hat{h}(t) = \begin{cases} c_1 e^{p_1 t} + c_2 e^{p_2 t} + \dots + c_N e^{p_N t} & t \geq 0 \\ 0 & t < 0 \end{cases} \quad (2-36)$$

(on suppose ici que les pôles p_i de $\hat{H}(s)$ sont simples).

$$\text{En posant } \hat{h}_i(t) = \begin{cases} c_i e^{p_i t} & t \geq 0 \\ 0 & t < 0 \end{cases} \quad (2-37)$$

la forme échantillonnée est alors :

$$h_i(k) = \begin{cases} T_e c_i e^{p_i T_e k} & k \geq 0 \\ 0 & k < 0 \end{cases} \quad (2-38)$$

dont la transformée en z est

$$H_i(z) = \frac{T_e c_i}{1 - e^{p_i T_e} z^{-1}} \quad (2-39)$$

$$\text{donc } H(z) = \sum_{i=1}^N H_i(z) = \sum_{i=1}^N \frac{T_e c_i}{1 - e^{p_i T_e} z^{-1}} \quad (2-40)$$

$$\text{et } H(e^{j\omega T_e}) = \sum_{i=-\infty}^{+\infty} \hat{H} \left[j\omega + jn \frac{2\pi}{T_e} \right] \quad (2-41)$$

où $H(e^{j\omega T_e})$ est la somme des versions de $\hat{H}(j\omega)$ répétées à des périodes de $2\pi/T_e$ (effet de l'échantillonnage de $\hat{h}(t)$). Pour éviter le phénomène du recouvrement de ces versions, cette méthode de synthèse se restreint seulement aux filtres à bandes étroites et à transitions raides[1].

b- Méthode de la transformation bilinéaire

Cette méthode consiste à construire pour chaque filtre analogique $\hat{H}(s)$, un filtre numérique $H(z)$ tel que:

$$H(z = e^{j\omega T_e}) = \hat{H}(s = \Omega(\omega)) \quad (2-42)$$

où $\Omega(\omega)$ est appelée la distorsion de fréquence.

Si ωT_e prend ses valeurs dans l'intervalle $[-\pi, \pi]$ (théorème

d'échantillonnage), alors $\Omega(\omega)$ varie de $]-\infty, +\infty [$. Cette transformation qui applique le plan complexe S dans le plan complexe Z , est appelée transformation bilinéaire; elle est définie par :

$$z = \frac{1+s}{1-s} \quad \Leftrightarrow \quad s = \frac{z-1}{z+1} \quad (2-43)$$

d'où
$$j\Omega(\omega) = \frac{e^{j\omega T_e} - 1}{e^{j\omega T_e} + 1}$$

c'est à dire
$$\omega_a = \Omega(\omega) = \text{tg}(\omega T_e / 2) \quad (2-44)$$

où ω_a est la fréquence analogique correspondant à la fréquence numérique ω .

Cette méthode évite le recouvrement et conserve la stabilité, cependant elle introduit une distorsion dans l'axe des fréquences qui est généralement remédiable par le procédé dit de préwarping [2].

2- Méthodes d'optimisation

a-Approximation de PADE

Pour synthétiser un filtre numérique dont la réponse fréquentielle est $H(\omega)$ avec $-\frac{\pi}{T_e} < \omega < \frac{\pi}{T_e}$, on l'approxime par un filtre RII causal :

$$\hat{H}(z) = \frac{\sum_{i=0}^M b(i) z^{-i}}{1 + \sum_{i=1}^N a(i) z^{-i}} = \frac{B_M(z)}{A_N(z)} \quad N \geq M \quad (2-45)$$

avec $a(0) = 1$

ou encore
$$\hat{H}(z) = \sum_{k=0}^{+\infty} \hat{h}(k) z^{-k}$$

La méthode consiste à égaliser $h(k)$ à $\hat{h}(k)$ pour $T = M + N + 1$ échantillons, c'est à dire pour $k = 0, 1, \dots, T$.

Donc avec
$$\hat{H}(z) \cdot A_N(z) = B_M(z)$$

on a
$$\sum_{i=0}^M a(i) h(k-i) = b(k) \quad \text{pour } k=0, 1, \dots, M$$
 (2-46)

$$\sum_{i=0}^N a(i) h(k-i) = 0 \quad \text{pour } k=M+1, \dots, N$$

Ainsi la résolution de ce système linéaire (2-46) permet de trouver les coefficients $a(i)$ et $b(i)$.

Cette méthode est fiable pour les filtres d'ordres élevés, mais en revanche, elle nécessite une capacité mémoire importante et la stabilité du filtre obtenu n'est pas toujours garantie [7].

b-Méthode des moindres carrés [6]

Cette méthode consiste à généraliser l'approximation précédente pour le reste des échantillons de la réponse impulsionnelle h (pour $k = T+1$ à ∞). On cherche les coefficients $a(i)$ et $b(i)$ avec une contrainte sur l'erreur d'approximation:

$$\begin{aligned} \varepsilon(z) &= (H(z) - \hat{H}(z)) U(z) \\ &= Y(z) - \hat{Y}(z) \end{aligned} \quad (2-47)$$

qui doit avoir une énergie minimale en supposant que l'entrée $u(k)$ soit un bruit blanc. Soit à minimiser :

$$W = \frac{1}{2\pi} \int_{-\pi}^{\pi} |\varepsilon(e^{j\theta})|^2 d\theta = E[e^2(k)] \quad (2-48)$$

avec $E[e^2(k)] = \sum_{k=0}^{+\infty} (h(k) - \hat{h}(k))^2$

Comme la minimisation de cette expression conduit à des équations non linéaires, on modifie la formulation de l'erreur de façon à rester dans le cas linéaire; c'est la méthode des moindres carrés modifiés.

Sachant que : $E(z) = \left[H(z) - \frac{B_M(z)}{A_N(z)} \right] U(z) \quad (2-49)$

on pose $\hat{E}(z) = E(z) A_N(z)$
 $= (H(z) A_N(z) - B_M(z)) U(z) \quad (2-50)$

donc il suffit de trouver les coefficients $a(i)$ et $b(i)$ qui minimisent l'énergie de l'erreur modifiée:

$$\begin{aligned} W &= \frac{1}{2\pi} \int_{-\pi}^{\pi} |H(e^{j\theta}) A_N(e^{j\theta}) - B_M(e^{j\theta})|^2 d\theta \\ &= \sum_{k=0}^{+\infty} (h(k) * a(k) - b(k))^2 \end{aligned} \quad (2-51)$$

La minimisation de W par rapport à $a(k)$ et $b(k)$ conduit au

système d'équation suivant:

$$\begin{matrix} \leftarrow & N+1 & \rightarrow & \leftarrow & M+1 & \rightarrow \\ \left[\begin{array}{ccc|ccc} & & & & & \\ & & & & & \\ & R & & - & H_0^T & \\ \hline & & & & & \\ & & & & & \\ & - & H_0 & & & I \\ \hline & & & & & \end{array} \right] \left[\begin{array}{c} A \\ \hline B \end{array} \right] = \left[\begin{array}{c} \alpha_n \\ 0 \\ \vdots \\ 0 \end{array} \right] \quad (2-52)$$

avec $[R]_{ij} = \sum h(i) h(i+|i-j|)$ $i, j = 0, \dots, N$

$[H_0]_{ij} = h(i-j)$ $i \leq j, i = 0, \dots, N$

$[A]_i = a(i-1)$ $i = 1, \dots, N+1$

$[B]_i = b(i-1)$ $i = 1, \dots, M+1$

et α_n l'énergie minimale de l'erreur.

La résolution de ce système linéaire peut être faite par n'importe quel algorithme de résolution des systèmes linéaires. Mais du point de vue temps de calcul et espace mémoire, l'algorithme de C.T.MULLIS et R.A.BOBERTS se révèle le plus efficace [7].

2-3 Transformation fréquentielle

C'est un outil mathématique très efficace dans la synthèse des différents types de filtres : passe-bas, passe-haut, passe bande, etc..., à partir d'un filtre prototype $H(z)$ généralement passe bas.

Le filtre construit est

$$G(z) = H(F(z)) \quad (2-53)$$

où $F(z)$ est appelée transformation fréquentielle et est donnée par

$$F(z) = \pm \prod_{i=1}^m \left[\frac{z - \alpha_i^*}{1 - \alpha_i z} \right] \quad (2-54)$$

avec $|\alpha_i| < 1 \quad \forall i$

Si $H(z)$ est d'ordre n , $G(z)$ sera d'ordre $m n$. $F(z)$ se comporte de la même manière que la variable z dans $H(z)$. Par conséquent la stabilité de $H(z)$ implique celle de $G(z)$.

A partir d'un filtre analogique prototype $\hat{H}(s)$, on peut synthétiser un filtre numérique $G(z)$ avec des caractéristiques bien spécifiées, par deux procédés: [1]

1- Soit appliquer une transformation fréquentielle analogique pour trouver $\hat{G}(s)$, ensuite utiliser la transformation bilinéaire pour obtenir $G(z)$.

2- Soit utiliser, d'abord, la transformation bilinéaire pour

trouver le filtre numérique prototype correspondant $H(z)$, ensuite, appliquer une transformation fréquentielle numérique pour obtenir $G(z)$.

Les transformations fréquentielles peuvent aussi être utiliser pour linéariser la réponse de phase d'un filtre numérique [2];[7]. Notons aussi que le spectre d'amplitude de telles transformations est constant sur toute la gamme fréquentielle, la raison pour laquelle $F(z)$ est appelée filtre passe tout.

2.4 Tableau comparatif entre les filtres RIF et RII

Filtres R I F	Filtres R I I
<p>* Ces filtres sont toujours stables, souvent utilisés dans les systèmes adaptatifs et multicanal [8].</p> <p>* leur formulation est simple et leur structure est non recursive.</p> <p>* Phase linéaire facilement obtenue.</p> <p>* Les erreurs de calculs dues aux effets des registres finis ne sont pas cumulatives et sont simples à analyser.</p> <p>* Temps de calcul long surtout si l'ordre est grand (bande de transition étroite).</p> <p>* Nécessitent un espace mémoire important si l'ordre est grand.</p>	<p>* La stabilité, dépend de la position des pôles par rapport au cercle unité[8].</p> <p>* Plus délicats au stade de leur concrétisation, sauf pour les filtres elliptiques[8].</p> <p>* La linéarité de la phase peut être seulement approchée à l'aide de systèmes d'égalisation du temps de propagation de groupe[8].</p> <p>* Les conséquences de la limitation de la longueur des registres sont graves car les erreurs de calcul sont cumulatives.</p> <p>* Calculs simples et moins longs même pour des transitions très raides.</p> <p>* Possibilité d'avoir des réponses fréquentielles très sélectives en jouant sur les pôles.</p> <p>* Nécessitent un faible espace mémoire.</p>

CHAPITRE III

**EFFETS DE LA LONGUEUR LIMITEE DES MOTS
SUR LE FILTRAGE NUMERIQUE**

CHAPITRE III

EFFETS DE LA LONGUEUR LIMITEE DES MOTS SUR LE FILTRAGE NUMERIQUE

3-1-Introduction

Dans les réalisations aussi bien matérielles que logicielles (simulation) des filtres numériques, on est amené à stocker les nombres dans des registres de longueurs finies. Par conséquent, les valeurs du signal et des paramètres caractérisant le filtre doivent être quantifiées avant d'être stockées. L'utilisation de registres limités oblige à accorder une grande attention aux effets qu'elle occasionne au filtrage; on peut en citer:

1-Le bruit de la conversion analogique-numérique.

2-Le bruit de calcul incorréolé.

3-L'inexactitude de la réponse du filtre due à la quantification des coefficients.

4-Le bruit de calcul corréolé ou les cycles limites.

Ces effets qui affectent la performance du filtre dépendent essentiellement:

-du type d'arithmétique utilisée dans l'algorithme du filtre

-du type de la quantification utilisée pour réduire les mots à la longueur désirée (par l'arrondi ou la troncature).

-de la structure choisie du filtre utilisé.

Indépendamment de l'encodage des signaux (c'est à dire, de l'arithmétique), les multiplications, et dans certains cas, les additions requièrent généralement une augmentation de la longueur du mot pour contenir le résultat de l'opération.

Comme le nombre d'opérations effectuées sur le signal reste fini, on peut accommoder l'accroissement de la longueur des mots par l'utilisation de registres plus longs que ceux du signal original afin de pouvoir stocker les résultats des opérations arithmétiques. Cependant, dans ces cas, de très longs registres peuvent être nécessaires et pour cette raison il est d'ordinaire en pratique, de réduire la taille des mots.

Pour les filtres non récurrents, une telle réduction de la taille des mots cause un signal d'erreur additif à la sortie similaire à celui dû à la quantification dans le convertisseur analogique-numérique.

Par contre, pour les filtres récurrents, une réduction de la longueur des mots est nécessaire à chaque boucle de retour, pour éviter une éventuelle augmentation de la longueur des mots. Cette opération non linéaire introduit des erreurs causant des effets graves [9].

D'autre part, les filtres utilisés le plus souvent en pratique sont ceux ayant de hauts facteurs de qualité donc des gains importants; en conséquence les signaux peuvent devenir plus grands que la valeur maximum représentable par le nombre de digits disponibles, en considérant la gamme dynamique désirée. Si un dépassement a lieu dans un filtre, le digit (généralement, binaire) le plus significatif doit être, à son tour, stocké dans un registre.

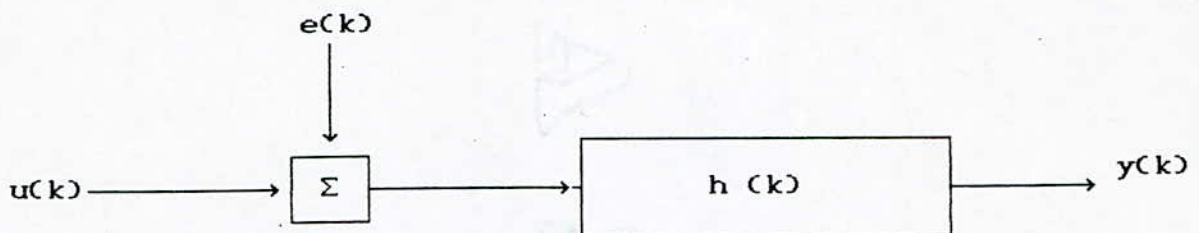
3-2 Conversion analogique numérique

Dans la plupart des applications de filtrage numérique, le signal original à filtrer est de nature analogique; c'est pourquoi un convertisseur analogique numérique (CAN) est toujours une partie intégrante du filtre numérique. Le CAN assure l'échantillonnage du signal continu avec une cadence au moins égale ou supérieure au double de la fréquence limite du signal et le codage des échantillons en une suite de nombres binaires (en général) de longueurs finies.

L'erreur entre un nombre réel $u(k)$ et sa représentation binaire finie; $[u(k)]_q$, est appelée bruit de quantification. Elle est supposée aléatoire et est définie par:

$$e(k) = u(k) - [u(k)]_q \quad (3-1)$$

On schématise l'effet de la quantification à l'entrée d'un filtre numérique de réponse impulsionnelle $h(k)$ par (Fig 3-1)



Fig(3-1)-Introduction de l'erreur de quantification

La sortie du filtre s'écrit:

$$y(k) = \underbrace{u(k) * h(k)}_{\substack{\text{signal de} \\ \text{sortie utile}}} + \underbrace{e(k) * h(k)}_{\text{effet du bruit}} \quad (3-2)$$

3-3 La quantification

La quantification est l'approximation de chaque valeur du signal par un multiple entier d'une quantité élémentaire q appelée pas de quantification. Cette opération consiste à faire passer le signal dans un organe ayant une caractéristique en marche d'escalier Fig(3-2).

Le centrage de cette caractéristique dépend de la manière avec laquelle l'approximation est faite:

-l'approximation par l'arrondi consiste à arrondir toute valeur comprise entre $(n - \frac{1}{2})q$ et $(n + \frac{1}{2})q$ à nq . Cette quantification rend la moyenne du signal d'erreur nulle (Fig(3-2-(a))).

-l'approximation par la troncature consiste à approcher par nq toute valeur comprise entre nq et $(n+1)q$ (Fig 3-2-(b)).

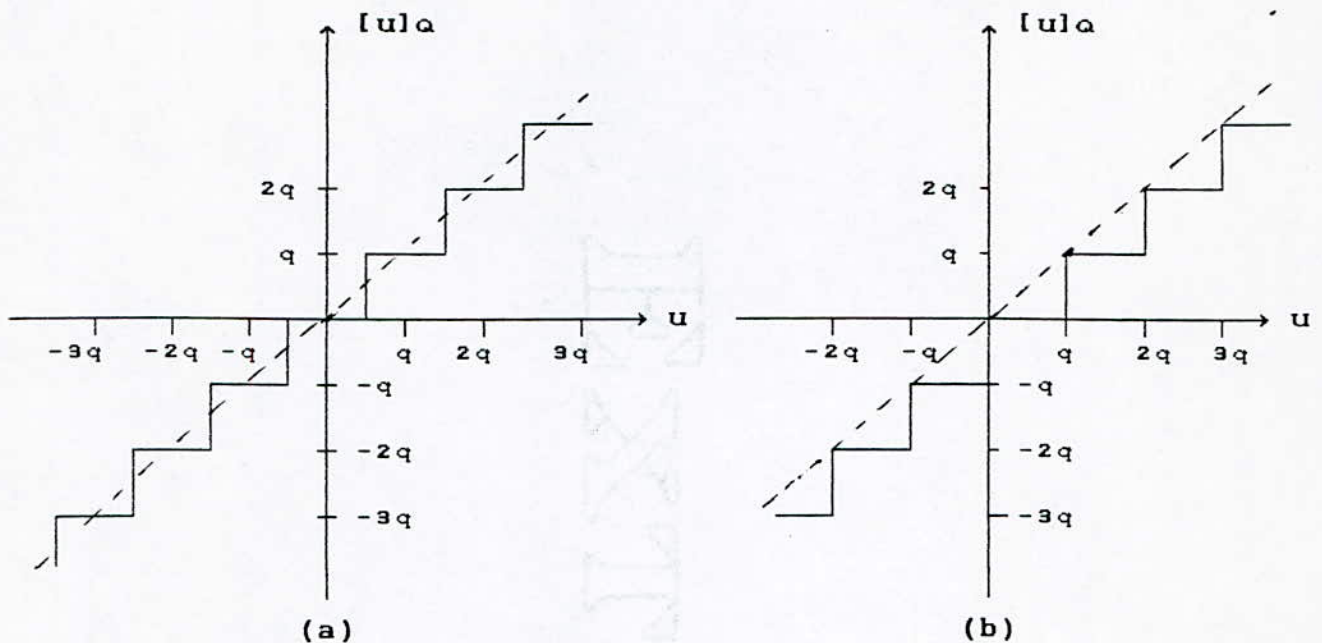


Fig (3-2) - Caractéristique du quantificateur
 (a) par l'arrondi, (b) par la troncature

Si la gamme des valeurs du signal couvre le domaine $[-\Delta, \Delta]$ et si l'on code (en binaire) chaque valeur par des mots de B bits,

$$q = \Delta 2^{-B+1} \quad (3-3)$$

dans ce cas la quantification est dite uniforme, (q constant)[8].

La distribution du bruit est uniforme, et sa puissance (variance) est estimée à:

$$P = \frac{q^2}{12} \quad (3-4)$$

3-4-Types de représentation arithmétique binaire

Il existe diverses façons d'établir la correspondance entre l'ensemble des amplitudes quantifiées et l'ensemble des nombres binaires qui doivent les représenter. Les signaux à coder ayant des valeurs, en général, positives et négatives, les représentations préférées sont celles qui conservent l'information du signe. Les plus courantes sont:

3-4-1 Représentation en virgule fixe

Généralement un nombre V s'exprime dans la base r par :

$$V = \sum_{i=-m}^n b_i r^i \quad \text{avec } 0 \leq b_i \leq r-1 \quad (3-5)$$

b_i entier

on le note dans cette base:

$$V = (b_n b_{n-1} \dots b_0 . b_{-1} \dots b_{-m}) \quad (3-6)$$

La représentation binaire correspond à $r = 2$ avec $b_i = 0$ ou 1 . Le point séparant V en deux parties s'appelle virgule binaire.

Dans la représentation, en virgule fixe, la virgule binaire se situe entre le premier et le second bits du registre les plus significatifs, où la première position représente le bit de signe.

Dans cette représentation un nombre V fractionnaire et signé peut s'écrire sous trois formes : [8], [10], [11]

1-Signe et valeur absolue:

$$V_{sva} = \begin{cases} 0.b_{-1}b_{-2}\dots b_{-m} & \text{si } V \geq 0 \\ 1.b_{-1}b_{-2}\dots b_{-m} & \text{si } V < 0 \end{cases} \quad (3-7)$$

2-Complément à 1 :

$$V_{c1} = \begin{cases} V & \text{si } V \geq 0 \\ 2 - 2^{-B} - |V| & \text{si } V < 0 \end{cases} \quad (3-8)$$

où B est la longueur du mot (nombre de bits)

3-Complément à 2 :

$$V_{c2} = \begin{cases} V & \text{si } V \geq 0 \\ 2 - |V| & \text{si } V < 0 \end{cases} \quad (3-9)$$

Le choix entre ces trois représentations arithmétiques est généralement dicté par les considérations des moyens matériels disponibles et de la programmation.

L'opération d'addition (ou soustraction) est directe dans le cas du complément à 2. L'arithmétique du signe et valeur absolue convient mieux à la multiplication car elle se fait simplement par la multiplication bit par bit des valeurs absolues et par l'ajustement du bit de signe du résultat.

Dans la plupart des réalisations de filtres numériques en virgule fixe, la position de la virgule binaire est supposée être à la droite du premier bit le plus significatif.

Ainsi la gamme des nombres représentables va de -1.0 à $1.0 - 2^{-(B-1)}$ où B est le nombre de bits du mot. C'est pourquoi souvent, les signaux sont normalisés afin de respecter la gamme choisie [10].

3-4-2 Représentation en virgule flottante

La représentation en virgule flottante d'un nombre consiste en deux parties; la mantisse et l'exposant, telle que un nombre V , dans cette représentation est noté :

$$V = m \times 2^e \quad (3-10)$$

où e est l'exposant (signé),

m est la mantisse signée et normalisée avec $\frac{1}{2} \leq m < 1$

Généralement le nombre de bits assigné à la mantisse est égal au trois-quart du nombre de bits total du mot [8].

Ce type de représentation permet une extension de la dynamique du fait de l'effet multiplicatif introduit par l'exposant, cependant il entraîne une complication des opérations arithmétiques et des circuits [11].

C'est pour cette raison qu'on se limite, dans la suite de ce travail, à l'étude du cas de la représentation en virgule fixe uniquement.

3-5 Erreurs de calcul et de dépassement

En plus de l'erreur causée par la conversion analogique-numérique à l'entrée du filtre qui est indépendante de la réalisation de celui-ci, des erreurs intrinsèques au filtre dues à la limitation de la longueur des mots (représentés en virgule fixe) introduisent des effets indésirables dans le filtrage. Ces effets causent trois types d'erreurs :

1-Erreurs de calcul dues à la quantification des produits

A chaque instant, dans un filtre numérique, un signal représenté par b_1 bits est multiplié par un coefficient représenté par b_2 bits en donnant, à la sortie du multiplieur, un résultat sur b_1+b_2

bits. Comme la longueur des registres est fixe pour tout le système, chaque sortie des multiplieurs doit être quantifiée avant la prochaine opération [10].

2- Erreurs de quantification des coefficients

Durant l'étape de synthèse, les coefficients de la fonction de transfert du filtre sont normalement évalués avec une grande précision. S'ils sont quantifiés, la réponse fréquentielle va différer considérablement de la réponse désirée et si, en plus, le pas de quantification est mal choisi, le filtre peut ne pas satisfaire les spécifications désirées.

Cet effet est appelé "sensibilité" des coefficients à la quantification c'est un effet déterministe qui peut être dominé par le bruit de calcul [1].

En effet la limitation du nombre de bits des coefficients se traduit par le fait qu'ils ne peuvent prendre qu'un nombre limité de valeurs, il s'en suit que les pôles ont un nombre limité de positions possibles à l'intérieur du cercle unité, il en est de même pour les zéros.

Ainsi la quantification à B bits de la valeur absolue des coefficients limité à 2^{2B} , le nombre de positions que peuvent prendre les pôles dans un quart du cercle unité et à 2^B le nombre de fréquences de coupure [8],[10].

3- Cycles limites dus aux effets non linéaires de la quantification :

Dans l'étude du bruit de calcul, on suppose toujours que les valeurs des échantillons du signal d'entrée du filtre numérique sont de même ordre de grandeurs que les différents multiples du pas de quantification q. Ce qui nous permet d'admettre que les échantillons du bruit de calcul sont statistiquement incorrélés aussi bien entre eux qu'avec la suite du signal d'entrée.

Cependant, le signal à l'entrée peut atteindre des valeurs faibles durant certains temps, par conséquent, les erreurs de quantification tendent à devenir fortement corrélés et peuvent causer l'instabilité du filtre par l'apparition d'une auto-oscillation appelée cycles limites [8],[10].

En effet, bien que les conditions de stabilité soient remplies et que le signal à l'entrée soit absent, un signal périodique peut

apparaître, généralement avec de faibles amplitudes. Ces oscillations tiennent au fait qu'en réalité le signal d'entrée n'est jamais nul en l'absence de données à l'entrée; le signal d'erreur dû à la quantification des nombres avant la mise en mémoire est appliqué au filtre.

Une borne peut être obtenue pour ces oscillations sachant que le signal d'erreur $e(k)$ possède en lui même, dans le cas d'arrondi (avec pas q), une borne donnée par :

$$|e(k)| \leq \frac{q}{2} \quad (3-11)$$

Si la réponse impulsionnelle du filtre est h , les auto-oscillations sont limitées par:

$$|y(k)| \leq \frac{q}{2} \sum_i |h(i)| \quad (3-12)$$

Cette borne est en fait très large; une estimation plus réaliste de l'amplitude des auto-oscillations est donnée par :

$$A_m = \frac{q}{2} \text{Max } |H(\omega)| \quad (3-12)$$

où $H(\omega)$ est la réponse impulsionnelle du filtre.

Pour que les cycles limites ne soient pas gênants, le nombre de bits des mémoires doit être suffisamment grand et le pas q suffisamment petit [1]. Il faut noter que ces cycles limites peuvent être éliminés par l'utilisation de la troncature au lieu de l'arrondi [10].

3-5-1 Oscillations de dépassement

Si l'amplitude d'un signal interne d'une réalisation en virgule fixe excède la gamme dynamique, un dépassement se produit et le signal de sortie va être distordu.

L'erreur causée par les dépassements est déterminée par la méthode utilisée pour représenter un nombre dépassant la gamme. Ce dernier peut être transformé par plusieurs façons, selon la caractéristique de dépassement choisie. Parmi lesquelles, on en cite:

a-La caractéristique de complément à 2 : utilisée souvent dans la représentation en complément à deux. Elle consiste à attribuer à chaque nombre dépassant la gamme, le complément à 2 du plus grand

(ou plus petit) nombre représentable.

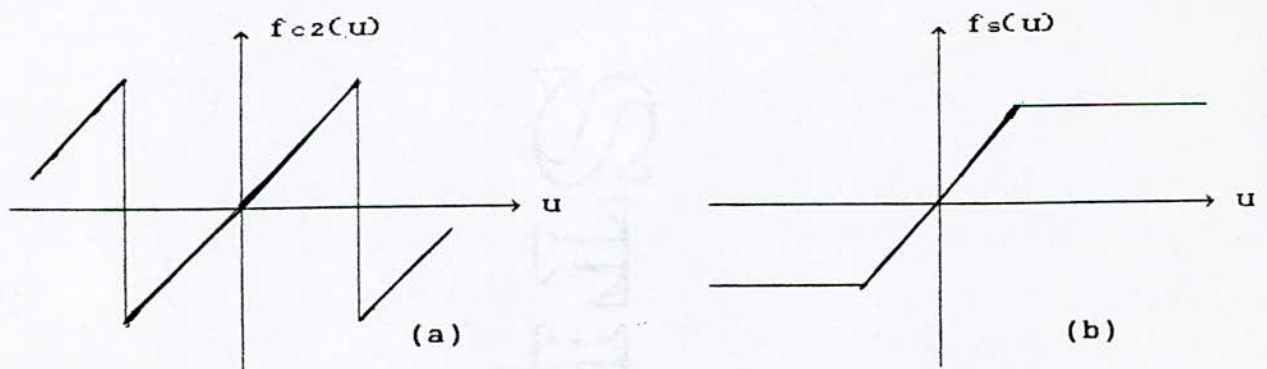
Cette caractéristique est périodique (voir Fig(3-3-(a)) et a pour propriété d'annuler les dépassements intermédiaires dans un accumulateur et d'obtenir des résultats de somme qui soient dans la gamme [1],[9].

b-La caractéristique de saturation qui consiste à remplacer le nombre dépassant la gamme par le plus grand (ou le plus petit) nombre représentable (Fig(3-3-(b))). L'erreur de dépassement de cette caractéristique est inférieure à celle du complément à 2, mais elle est difficile à réaliser matériellement [1],[9].

Lorsqu'un dépassement se produit à l'entrée, la propagation de l'erreur peut causer d'autres dépassements. Aussi, après un dépassement interne la sortie du filtre peut, en dépendant des pôles de celui-ci, devenir indépendante du signal d'entrée; c'est l'auto-oscillation de dépassement. Les oscillations de dépassement sont causées par la non linéarité de la caractéristique de dépassement utilisée [1],[9],[10].

Les dépassements peuvent être minimisés en élargissant la gamme des représentations possibles; et ceci en augmentant la valeur de Δ de (3-3).

Cependant, si le nombre de niveaux de quantification reste le même la valeur du pas q doit être augmentée aussi; en conséquence, une diminution de la possibilité de dépassement n'est obtenue qu'au prix d'une augmentation de l'erreur de quantification.



Fig(3-3) Caractéristiques de dépassement :
(a) complément à 2, (b) saturation

CHAPITRE IV

**REPRESENTATION DES FILTRES NUMERIQUES
PAR LES VARIABLES D'ETAT**

CHAPITRE IV

REPRESENTATION DES FILTRES NUMERIQUES PAR LES VARIABLES D'ETAT

4-1- Introduction

Jusqu'à présent, le filtre numérique a été représenté par sa caractéristique d'entrée-sortie, c'est à dire sa fonction de transfert ou sa réponse impulsionnelle qui caractérisent complètement les propriétés et les spécifications du filtre à réaliser. Mais étant donné que tout système numérique n'effectue que trois types d'opérations: les multiplications, les additions et les décalages; la réalisation du filtre numérique exige en plus de la détermination de la fonction de transfert, la connaissance de l'ordre des opérations élémentaires à exécuter. Donc il devient impératif de concevoir une représentation plus détaillée du filtre afin de décrire exactement l'opération du filtrage.

La connaissance de la manière dont s'effectuent les opérations à l'intérieur du filtre permet de mieux analyser les erreurs dues à l'utilisation de registres de longueurs finies et par conséquent de trouver répondant aux spécifications du filtre et présentant un meilleur rapport signal/bruit.

Pour ce fait, on utilise la représentation par les variables d'état pour décrire la réalisation du filtre numérique, qui constitue en plus d'un outil mathématique puissant, un moyen d'analyse souple et efficace que nous allons découvrir par la suite.

4-2 Représentation d'état des filtres numériques

Tout système linéaire d'ordre N , à l'instant k , auquel est appliquée une suite $u(k)$ à l'entrée et qui fournit une suite $y(k)$ en sortie, est défini par le couple d'équation dites d'état suivant:

$$\begin{aligned} x(k+1) &= A x(k) + B u(k) \\ y(k) &= C x(k) + D u(k) \end{aligned} \quad (4-1)$$

où $x(k)$ est un vecteur $N \times 1$ dont chaque élément est une variable d'état du système (ou filtre).

A est une matrice $N \times N$ appelée "matrice du système" (filtre).

B est un vecteur $N \times 1$, c'est le vecteur de "commande".

C est un vecteur $1 \times N$ appelée vecteur "d'observation".

D est le coefficient de la transition directe de l'entrée à la sortie.

L'état du filtre, à l'instant k , en fonction de l'état initial $x(0)$ est

$$x(k) = A^k x(0) + \sum_{i=1}^k A^{i-1} B u(k-i) \quad (4-2)$$

La fonction de transfert du système est obtenue en prenant la transformée en Z des équations (4-1):

$$\begin{aligned} (z I - A)^{-1} X(z) &= B U(z) \\ Y(z) &= C X(z) + D U(z) \end{aligned} \quad (4-3)$$

En éliminant $X(z)$, on obtient:

$$H(z) = C (z I - A)^{-1} B + D \quad (4-4)$$

Les pôles de la fonction de transfert ainsi obtenue, ne sont en fait que les racines du polynôme caractéristique de la matrice A tel que:

$$a(z) = \det(z I - A) \quad (4-5)$$

Ceux sont donc les valeurs propres de la matrice A.

Pour assurer la stabilité du système, il suffit que les modules des valeurs propres de A soient inférieurs à l'unité. Ainsi, en régime libre (entrée nulle):

$$x(k) = A^k x(0) \quad (4-6)$$

tendra vers zéro quand k tend vers l'infini. En régime établi, l'équation (4-2) devient pour les systèmes stables:

$$x(k) = \sum_{i=1}^{+\infty} A^{i-1} B u(k-i) \quad (4-7)$$

(quand k tend vers l'infini)

La réponse impulsionnelle $h(k)$ du filtre peut être exprimée

en fonction des paramètres d'états; sachant que la sortie du filtre s'écrit:

$$y(k) = \sum_{l=0}^{+\infty} h(l) u(k-l) \quad (4-8)$$

et par ailleurs ,

$$\begin{aligned} y(k) &= C x(k) + D u(k) \\ &= C \sum_{i=1}^{+\infty} (A^{i-1} B u(k-i)) + D u(k) \\ &= \sum_{i=1}^{+\infty} C A^{i-1} B u(k-i) + D u(k) \end{aligned} \quad (4-9)$$

Enfin, par identification des équations (4-8) et (4-9), on obtient

$$h(k) = \begin{cases} 0 & , k < 0 \\ D & , k = 0 \\ C A^{k-1} B & , k > 0 \end{cases} \quad (4-10)$$

* Changement de coordonnées

La fonction de transfert (ou la réponse impulsionnelle) spécifie la relation entrée-sortie du filtre sans pour autant apporter une information sur sa structure interne.

Comme le vecteur d'état $x(k)$ représente N variables d'état internes, on peut changer le système de coordonnées de ce vecteur dans l'espace d'état sans changer la relation d'entrée-sortie (ou la réponse impulsionnelle du filtre). Dans ce sens, il existe une certaine latitude dans le choix des paramètres d'état pour lesquels les valeurs propres de la matrice A imposées par le filtre restent inchangées.

En appliquant une transformation T , une matrice non singulière $N \times N$, le nouveau vecteur d'état est donnée par:

$$x'(k) = T^{-1} x(k) \quad (4-11)$$

Les paramètres d'état se transforment alors comme suit

$$(A, B, C, D) \longrightarrow (A' = T^{-1} A T, B' = T^{-1} B, C' = C T, D) \quad (4-12)$$

On peut vérifier que ce type de transformation conserve la fonction de transfert et donc la réponse impulsionnelle [1].

Une structure canonique ou directe du filtre est définie par des paramètres de la forme:

$$A = \begin{bmatrix} 0 & & & & \\ \vdots & & & & \\ & & I & & \\ & & & & \\ -\frac{a_N}{a_0} & -\frac{a_{N-1}}{a_0} & \dots & \dots & -\frac{a_1}{a_0} \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \quad C^T = \begin{bmatrix} b_N - \frac{b_0}{a_0} & a_1 \\ \vdots & \\ b_1 - \frac{b_0}{a_0} & a_0 \end{bmatrix}, \quad D = \frac{b_0}{a_0} \quad (4-13)$$

où a_i et b_i ($i=0, \dots, N$) sont les coefficients du dénominateur de la fonction de transfert du filtre.

4-3 Matrice K et la stabilité de LYAPUNOV

La représentation des filtres numériques linéaires par les variables d'état est plus efficace quand il s'agit de calculer les quantités qui dépendent de la structure interne du filtre et de connaître la manière dont ces grandeurs changent quand la structure est modifiée.

Ces quantités à calculer sont généralement statistiques et interviennent quand les entrées appliquées aux filtres sont stationnaire au sens large [1].

Parmi ces quantités, on s'intéresse surtout à la matrice de covariance de l'état du filtre décrit par les paramètres (A,B,C,D) et qui est donnée par :

$$K = E [x(k) x^T(k)] \quad (4-14)$$

En supposant que la stabilité du système fait que

$$x(k+1) = x(k) \quad \text{pour } k \rightarrow +\infty \quad (4-15)$$

donc les éléments de K sont :

$$K_{ij} = E [x_i(k) x_j(k)] \quad (4-16)$$

En remplaçant l'équation (4-7) dans (4-14), on a:

$$K = \sum_{l=0}^{\infty} \sum_{m=0}^{\infty} (A^l B) r_{uu}(l-m) (A^m B)^T \quad (4-17)$$

avec $r_{uu}(\tau) = E [u(k) u(k+\tau)]$

C'est la fonction d'autocorrélation de l'entrée.

Pour une entrée bruit blanc centré et normalisé, on a

$$r_{uu}(\tau) = \delta(\tau) \quad (4-18)$$

La double somme de l'équation (4-17) devient alors:

$$K = \sum_{i=0}^{+\infty} (A^i B) (A^i B)^T \quad (4-19)$$

On peut éviter les sommes infinies dans le calcul de la matrice K en utilisant les équations d'état:

$$\begin{aligned} K &= E [x(k+1) x^T(k+1)] \\ &= E [(A x(k) + B u(k))(A x(k) + B u(k))^T] \quad (4-20) \end{aligned}$$

Et sachant que $x(k)$ et $u(k)$ sont incorrélés pour des entrées bruit blanc, alors (4-20) devient:

$$K = E [A x(k) x^T(k) A^T + B u(k) u^T(k) B^T]$$

d'où

$$K = A K A^T + B B^T \quad (4-21)$$

Cette dernière expression est appelée equation de LYAPUNOV.

Un cas particulier de la théorie de la stabilité de LYAPUNOV est:

Si le couple (A,B) est contrôlable, c'est à dire que matrice de contrôlabilité

$$M = \begin{bmatrix} B, AB, A^2 B, \dots, A^{N-1} B \end{bmatrix} \quad (4-22)$$

est de rang N, et si la matrice K est définie positive alors la matrice A doit être stable [1].

4-4- Structures éliminant les oscillations de dépassement

Les oscillations de dépassement se produisent si des vecteurs d'état augmentent en amplitude après une multiplication par la matrice A du filtre. Comme la caractéristique de dépassement diminue l'amplitude du vecteur d'état si celle-ci dépasse la gamme, des oscillations répétées ne peuvent se produire à moins que le processus de la multiplication par A augmente la grandeur du vecteur d'état.

Une condition nécessaire pour la production des oscillations de

dépassement est que la grandeur de $A \cdot x$ soit supérieure à celle de x , pour un x donné de l'espace d'état.

on peut trouver une autre structure d'état caractérisée par (A, B, C, D) pour laquelle la grandeur de x n'augmente pas du fait de la multiplication par A . On définit la norme d'un vecteur comme étant

$$\|x\| = (x^T d x)^{1/2} \quad (4-23)$$

où d est une matrice diagonale positive.

La norme de la matrice A est définie par:

$$\begin{aligned} \|A\| &= \max_{x \neq 0} \frac{\|A x\|}{\|x\|} \\ &= \max_{x \neq 0} \left(\frac{x^T A^T d A x}{x^T x} \right)^{1/2} \end{aligned} \quad (4-24)$$

Pour éviter les oscillations de dépassement pour une entrée nulle (ou constante), il suffit que:

$$\|A x\| \leq r \|x\| \quad \text{avec} \quad 0 < r \leq 1 \quad (4-25)$$

Autrement dit, il suffit de trouver la matrice d telle que [1]:

$$Q = r^2 d - A^T d A \quad (4-26)$$

soit une matrice définie positive.

Parmi les structures qui satisfont la condition : Q définie positive; les structures normales ($A A^T = A^T A$), de gain de bruit de calcul minimal, et en treillis [13]. Les formes directes (canoniques) sont susceptibles de produire des oscillations de dépassement qui dépendent de la position des pôles.

Dans le cas des filtres du second ordre, si la matrice A (2×2) dont les valeurs propres λ vérifient: $|\lambda| < 1$, alors il existe une matrice diagonale définie positive d pour laquelle Q de (4-26) est définie positive, si et seulement si, les conditions suivantes sont satisfaites:

$$\begin{aligned} & a_{12} a_{21} \geq 0, \text{ ou} \\ \text{si } & a_{12} a_{21} < 0 \text{ alors } |a_{11} - a_{22}| + |\det(A)| < 1 \end{aligned} \quad (4-27)$$

tion de dépassement [14], [15].

X 4-5- Normalisation des filtres numériques en virgule fixe

Pour éviter les dépassements dans les registres internes qui causent des erreurs importantes, on normalise convenablement la réalisation du filtre.

Normaliser revient à mettre toutes les valeurs numériques des variables internes dans une gamme appropriée à la réalisation matérielle. La gamme d'une variable interne est nécessairement limitée par le fait de l'utilisation des registres de longueur finie.

4-5-1- Règles de normalisation

Dans la représentation en virgule fixe, une variable interne $v(k)$ telle que

$$v(k) = (f * u)(k) \quad (4-28)$$

avec $f(k)$ la réponse impulsionnelle entre l'entrée $u(k)$ et la variable d'état $v(k)$;

est bornée par le fait que l'entrée $u(k)$ est limitée par :

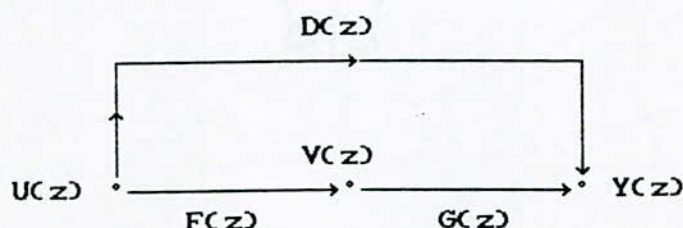
$$|u(k)| \leq \Delta \quad (4-29)$$

où Δ dépend du pas de quantification (d'après l'équation (3-3)) et qui est généralement normalisé à 1 (dans la représentation en virgule fixe).

La gamme des valeurs de v dépend ainsi de la nature de l'entrée $u(k)$ et de la suite $f(k)$. En effet, si un filtre caractérisé par sa fonction de transfert

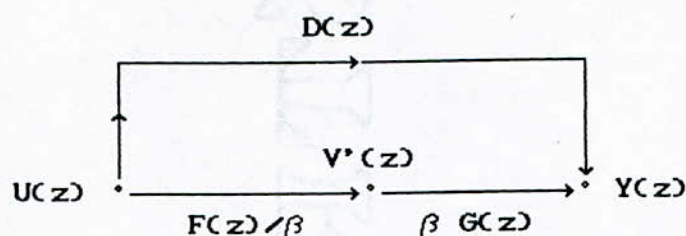
$$H(z) = D(z) + F(z) G(z) \quad (4-30)$$

et son graphe de fluence (Fig(4-1)):



Fig(4-1) Filtre non normalisé

Il est alors normalisé (c'est à dire que $v'(k)$ devient dans la gamme), de la façon suivante (Fig(4-2))



Fig(4-2) Filtre normalisé

De cette manière, $H(z)$ reste invariant et

$$|v'(k)| \leq \frac{1}{\beta} \sum_m |f(m)| |u(k-m)| \quad (4-31)$$

et comme $|u(k)| \leq 1$ alors $|v'(k)| \leq \frac{1}{\beta} \sum_m |f(m)| \quad (4-32)$

Donc β est choisi de façon à vérifier la relation (4-32) on peut choisir la norme 1 telle que:

$$\beta = \|f\|_1 = \sum_{l=0}^{\infty} |f(l)| \quad (4-33)$$

Mais étant donnée que cette norme constitue une borne largement conservatrice de la gamme et que la norme 2:

$$\|f\|_2 = \left[\sum_{l=0}^{\infty} f^2(l) \right]^{1/2} \leq \|f\|_1 \quad (4-34)$$

On peut donc choisir

$$\beta = \delta \|f\|_2 \quad (4-35)$$

Avec δ un paramètre choisi subjectivement afin d'obtenir une bonne représentation des valeurs de la variable $v'(k)$ dans la gamme, d'éviter ainsi les oscillations dues aux amplitudes faibles (ou cycles limites), et de réduire les erreurs d'arrondi.

En conclusion, un filtre est normalisé s'il vérifie la contrainte de normalisation définie par une des règles de normalisation dont on cite:[1]

1- Normalisation par la norme 1:

$$\|f\|_1 = 1 \quad (4-36)$$

2- Normalisation par la norme 2:

$$\delta \|f\|_2 = 1 \quad (4-37)$$

Dans notre cas on opte pour la norme 2 qui permet à la fois de

conserver la gamme et de réduire la probabilité des cycles limites et les erreurs d'arrondi par le choix d'une valeur adéquate de δ .

D'après les résultats de simulation, $\delta = 4$ constitue un compromis optimal entre les erreurs d'arrondi et les erreurs de dépassement.

Mais généralement on prend $\delta = 1$.

4-5-2 Normalisation des paramètres d'état

La normalisation d'un filtre par la norme 2, consiste à normaliser le vecteur d'état qui s'exprime par :

$$x(k) = \sum_{l=0}^{\infty} A^l B u(k-1-l) \quad (4-38)$$

La réponse impulsionnelle de l'entrée à $x(k)$ est alors :

$$f(k) = \begin{cases} 0 & , k \leq 0 \\ A^{k-1} B & , k > 0 \end{cases} \quad (4-39)$$

et donc $F(z) = (zI - A)^{-1} B$

$f(k)$ est un vecteur $N \times 1$ à l'aide duquel on peut exprimer la matrice K de (4-20) pour une entrée bruit blanc normalisé comme suit :

$$K = \sum_{k=0}^{\infty} f(k) f^T(k) \quad (4-40)$$

avec $K_{ij} = \langle f_i, f_j \rangle = \sum_{k=0}^{\infty} f_i(k) f_j(k) \quad (4-41)$

donc $K_{ii} = \sum_{k=0}^{\infty} f_i^2(k) = \|f\|_2^2 \quad (4-42)$

Donc pour normaliser ce filtre, il faut lui appliquer une transformation T diagonale pour laquelle la contrainte de normalisation (norme 2) soit :

$$\delta \sqrt{K'_{ii}} = 1 \quad (4-43)$$

dans la nouvelle représentation d'état.

Donc il suffit de prendre :

$$T_{ii} = \delta \sqrt{K_{ii}} \quad (4-44)$$

car d'après les équations (4-12) et (4-20), on a :

$$K \longrightarrow K' = T^{-1} K T^{-T} \quad (4-45)$$

et donc

$$K'_{ij} = \frac{K_{ij}}{T_{ii} T_{jj}}$$

d'où

$$K'_{ii} = 1/\delta^2 \quad (4-46)$$

La transformation de normalisation T doit être appliquée aussi aux paramètres d'états (A, B, C, D) suivant la relation (4-12).

En résumé, pour normaliser un filtre numérique:

- i)- on résout l'équation de LYAPUNOV (4-20) pour trouver K .
- ii)- on utilise la relation (4-44) pour construire T .
- iii)- on applique T aux paramètres (A, B, C, D) .

On remarque que la normalisation change la représentation d'états du filtre, donc sa structure interne et par conséquent le bruit de calcul.

4-6 Analyse du bruit de calcul dans le filtre numérique RII

Le bruit de calcul est l'erreur due, essentiellement, à la quantification des résultats des multiplications au niveau des accumulateurs durant l'opération de filtrage.

On assimile le quantificateur (dans le cas d'arrondi) à une source de bruit blanc centré de variance $q^2/12$, (où q est le pas de quantification), et de distribution uniforme sur $[-q/2, q/2]$. Ce modèle de bruit remplit deux conditions:

- a) Les sources des différents accumulateurs sont incorrélatées.
- b) Chaque source est incorrélatée avec l'entrée.

On suppose dans cette analyse que la probabilité des dépassements est très faible après une normalisation appropriée des registres internes du filtre.

Ainsi le graphe de fluence du filtre est donnée par Fig(4-3)

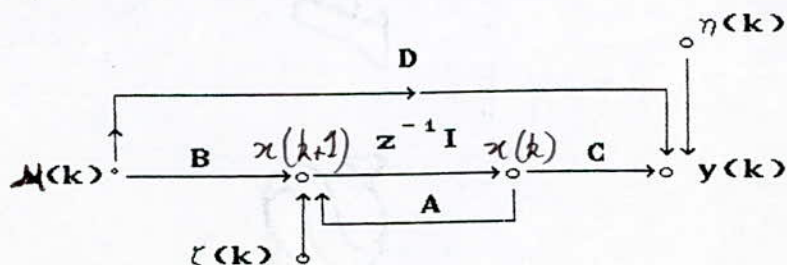


Fig (4-3) Graphe de fluence d'un filtre numérique RII

Où les sources d'erreurs de quantification sont associées aux

noeuds internes et au noeud de la sortie.

Les equations d'etat du filtre deviennent donc :

$$\begin{aligned} x(k+1) &= A x(k) + B u(k) + \eta(k) \\ y(k) &= C x(k) + D u(k) + \zeta(k) \end{aligned} \quad (4-47)$$

où $\zeta(k)$ est la source de bruit correspondant à la quantification des produits $A x(k)$ et $B u(k)$, et $\eta(k)$ est celle associée aux produits $Cx(k)$ et $Du(k)$. ($\eta(k)$ est un scalaire).

Comme la normalisation concerne la fonction de transfert entrée-variante d'etat $F(z)$ qu'on schématise par :

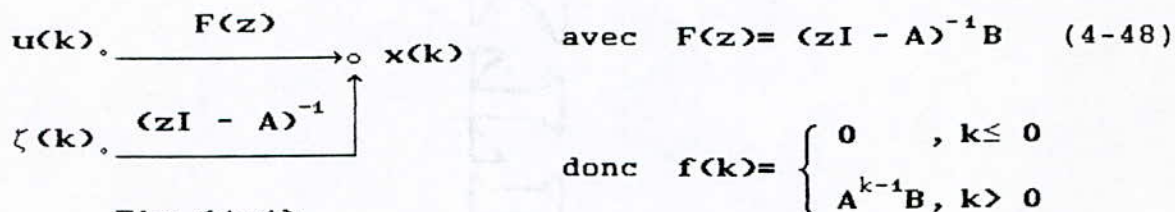


Fig (4-4)

Remarque

On considère que la contribution de $\zeta(k)$ est négligeable dans la gamme de $x(k)$ où les erreurs de dépassement sont prépondérantes pour un filtre non normalisé.

Le bruit de calcul concerne la fonction de transfert variable d'etat-sortie $G(z)$. Le principe de superposition permet de représenter $G(z)$ comme suit :

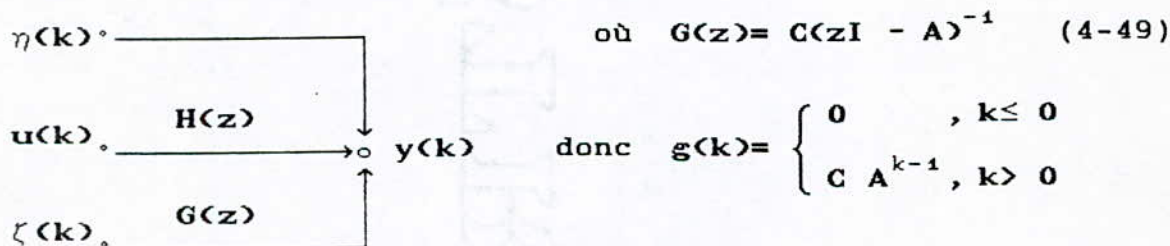


Fig (4-5)

$g(k)$ représente la réponse impulsionnelle du système variable d'etat-sortie. D'après les deux equations d'etat (4-1) et (4-47) l'erreur à la sortie est :

$$\Delta y(k) = \sum_{i=0}^{k-1} g(k-i) \zeta(i) + \eta(k)$$

$$= \sum_{i=0}^{k-1} C A^{k-i-1} \zeta(i) + \eta(k) \quad (4-50)$$

Donc la variance de cette erreur s'exprime comme:

$$\begin{aligned} E \left[\Delta y^2(k) \right] &= E \left[\left\{ \sum_{i=0}^{k-1} C A^{k-i-1} \zeta(i) + \eta(k) \right\} \right. \\ &\quad \left. \left\{ \sum_{j=0}^{k-1} C A^{k-j-1} \zeta(j) + \eta(k) \right\}^T \right] \\ &= \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} \left[C A^{k-i-1} \right] E \left[\zeta(k) \zeta^T(k) \right] \left[C A^{k-j-1} \right]^T \\ &\quad + E \left[\eta^2(k) \right] \end{aligned} \quad (4-51)$$

Comme les différentes composantes de $\zeta(k)$ sont incorrélées entre elles, on écrit:

$$E \left[\zeta(k) \zeta^T(k) \right] = \frac{q^2}{12} Q \quad (4-52)$$

où $Q = \text{diag} (\nu_1 \nu_2 \dots \nu_N)$

ν_i est le nombre de sources de bruit (ou de multiplieurs) associées au noeud de la i ème variable d'état. Dans ce qui suit, on notera ν_s , le nombre de sources associées au noeud de la sortie.

La variance de l'erreur de sortie à la limite ($k \rightarrow \infty$), s'écrit :

$$\sigma_y^2 = E \left[\Delta y^2(k) \right] = \frac{q^2}{12} \sum_{l=0}^{\infty} (C A^l) Q (C A^l)^T + \frac{q^2}{12} \nu_s \quad (4-53)$$

En notant par Tr la trace d'une matrice, (4-53) s'écrira:

$$\begin{aligned} \sigma_y^2 &= \frac{q^2}{12} \left[\sum_{l=0}^{\infty} \text{Tr} \left[(C A^l)^T (C A^l) Q \right] + \nu_s \right] \\ &= \frac{q^2}{12} \left[\text{Tr}(W Q) + \nu_s \right] \end{aligned} \quad (4-54)$$

$$\text{où } W = \sum_{l=0}^{\infty} (C A^l)^T (C A^l) = \sum_{k=1}^{\infty} g^T(k) g(k) \quad (4-55)$$

En procédant de la même manière que pour la matrice K (voir section 4-3), on élimine la somme infinie dans (4-55) et on obtient:

$$W = A^T W A + C^T C \quad (4-56)$$

L'équation (4-54) peut aussi s'écrire:

$$\begin{aligned}
 \sigma_y^2 &= \frac{q^2}{12} \left[\sum_{i=1}^N \nu_{ii} W_{ii} + \nu_s \right] \\
 &= \frac{q^2}{12} \left[\sum_{j=1}^N \nu_{ii} \left[\sum_{k=1}^{\infty} \varepsilon^2(k) \right] + \nu_s \right] \\
 &= \frac{q^2}{12} \left[\sum_{j=1}^N \nu_{ii} \|\varepsilon_j\|_2^2 + \nu_s \right] \quad (4-57)
 \end{aligned}$$

Dans le cas le plus défavorable pour un filtre d'ordre N, le nombre des multiplieurs au niveau de chaque noeud est au maximum :

$$\nu_s = \nu_i = N+1 \quad \text{pour tout } i$$

Dans ce cas, la variance est:

$$\sigma_y^2 = \frac{q^2}{12} (N+1) \left[\text{Tr}(W) + 1 \right] \quad (4-58)$$

Et pour le cas le plus favorable, on utilise des registres de longueur double pour accumuler les produits, donc une seule quantification est effectuée au niveau de chaque noeud, donc l'expression de la variance devient:

$$\sigma_y^2 = \frac{q^2}{12} \left[\text{Tr}(W) + 1 \right] \quad (4-59)$$

Dans tous les cas, on définit le gain de bruit comme étant le terme

$$G = \text{Tr}(W) = \sum_{i=1}^N W_{ii} \quad (4-60)$$

En effectuant un changement de coordonnées dans l'espace d'état par une transformation T, non singulière, N x N; la matrice W se transforme d'après la relation (4-56) en:

$$W' = T^T W T \quad (4-61)$$

$$\text{Le gain de bruit devient alors : } G' = \text{Tr}(T^T W T) \quad (4-62)$$

En particulier, si T est la transformation de normalisation définie précédemment par l'équation (4-44), W' aura pour composantes :

$$W'_{ij} = [T^T W T]_{ij} = \delta^2 \sqrt{K_{ii} K_{jj}} W_{ij} \quad (4-63)$$

Par conséquent, la variance de la sortie devient:

$$\sigma_y^2 = \frac{q^2}{12} (N + 1) \left[\delta^2 \sum_{i=1}^N K_{ii} W_{ii} + 1 \right] \quad (4-63)$$

Cette expression donne la variance de l'erreur de sortie en fonction des paramètres du filtre non normalisé.

L'objectif essentiel de tout ce qui va suivre (voir chapitre V) sera la minimisation de cette variance et plus précisément la quantité $\sum_{i=1}^N K_{ii} W_{ii}$, simplement par une transformation de coordonnées tout en conservant les contraintes de normalisation.

4-7 Remarque importante:

Si la valeur de δ est importante, la probabilité de dépassement dans la gamme des variables internes est faible, cependant, l'erreur de calcul à la sortie, d'après l'équation (4-64) sera importante à cause de la mauvaise correspondance entre la gamme dynamique et les nombres à représenter.

Si la valeur de δ est faible, l'erreur due au dépassement va prédominer dans la variance de l'erreur en sortie du filtre.

Il a été montré [16], par simulation, que la valeur $\delta = 4$ constitue un compromis optimal entre les erreurs de dépassement et les erreurs de quantification des opérations arithmétiques.

CHAPITRE V

STRUCTURES OPTIMALES

STRUCTURES OPTIMALES

5-1- Minimisation du gain de bruit dans le cas de registres de longueurs égales

La minimisation du gain de bruit est relative à la variance de l'erreur de calcul à la sortie du filtre. Elle n'aura de sens que si l'on suppose que le facteur de normalisation a été convenablement choisi afin de rendre les dépassements peu probables et l'augmentation du bruit de calcul qui en découle négligeable.

Minimiser le gain de bruit d'un filtre numérique décrit par ses paramètres d'états (A, B, C, D) et ses matrices K et W, revient à trouver une structure par une transformation T, qui change les paramètres du filtre en $(T^{-1}AT, T^{-1}B, CT, D)$ et K et W en $T^{-1}K T^{-T}$ et $T^T W T$, respectivement, de sorte que le gain de bruit soit minimal avec la contrainte de normalisation.

Comme la quantité $\sum_{i=1}^N K_{ii} W_{ii}$ est invariante à une transformation diagonale, il est donc possible de minimiser d'abord le gain de bruit ensuite d'appliquer une transformation de normalisation à la structure trouvée.

Actuellement, il existe deux méthodes de base permettant de construire les structures à gain de bruit minimal, la méthode de Mullis-Roberts [16] et la méthode de Hwang [17]. Bien que ces deux méthodes sont différentes du point de vue du traitement mathématique du problème, elles aboutissent chacune à un même gain minimal.

A - Methode de Mullis-Roberts [16]

Cette méthode traite deux cas; le cas des registres de longueurs égales et celui des registres de longueurs optimales. On se limitera à l'étude du premier cas uniquement.

La méthode consiste à rechercher le gain de bruit optimale à partir des conditions de son existence qui permettent d'établir la transformation réalisant la structure optimale.

Si le pas de quantification $q = 1.2^{-B+1}$ et B le nombre de bits des registres considérés, l'expression de la variance en sortie

(4-64) devient :

$$\sigma_y^2 = \frac{N(N+1)}{3} \left(\frac{\delta}{2^B} \right)^2 \left[\frac{1}{N} \sum_{i=1}^N K_{ii} W_{ii} \right] + \frac{N+1}{3 \cdot 2^{2B}} \quad (5-1)$$

Si N, δ et B sont donnés alors il suffit de minimiser le gain de bruit donné par:

$$G = \sum_{i=1}^N K_{ii} W_{ii}$$

A cet effet, on utilise la propriété des matrices K et W : dans le cas où celles-ci sont réelles, symétrique et définies positives, l'inégalité suivante est vérifiée :

$$\left[\frac{1}{N} \sum_{i=1}^N K_{ii} W_{ii} \right] \geq M_\alpha^2 \quad (5-2)$$

où
$$M_\alpha = \frac{1}{N} \sum_{i=1}^N \mu_i \quad (5-3)$$

et $\mu_1^2, \mu_2^2, \dots, \mu_N^2$ sont les valeurs propres de la matrice produit KW , et dont les racines carrées sont appelées les modes de second ordre [1].

Pour que la borne inférieure soit atteinte, il est nécessaire et suffisant de satisfaire les deux conditions suivantes [1],[16]:

- 1) $D_0^{-1} K D_0^{-1} = D_0 W D_0$ où D_0 est une matrice diagonale autrement dit $K = D' W D'$ telle que $D' = D_0^2$
- 2) $K_{ii} W_{ii} = K_{jj} W_{jj} \quad \forall i, j = 1, \dots, N$

Si K est normalisée ($K_{ii} = 1/\delta^2$) alors $W_{ii} = W_{jj} \quad \forall i, j = 1, \dots, N$

La condition 1) implique que si $K' = D_0^{-1} W D_0^{-1} = W'$, alors [6]

$$\text{Tr}(K') = \text{Tr}(W') = \sum_{i=1}^N \mu_i = N M_\alpha \quad (5-4)$$

La condition 2) implique l'égalité suivante:

$$K'_{ii} W'_{ii} = K_{ii} W_{ii} = K_{jj} W_{jj} = \text{constante} = \alpha \quad \forall i, j = 1, \dots, N \quad (5-5)$$

donc
$$\text{Tr}(K') = N M_\alpha = N \sqrt{\alpha} \quad \text{et} \quad \alpha = M_\alpha^2 \quad (5-6)$$

d'où
$$\sum_{i=1}^N K'_{ii} W'_{ii} = \alpha N = N M_\alpha^2$$

donc on a bien l'égalité
$$\frac{1}{N} \sum_{i=1}^N K'_{ii} W'_{ii} = M_\alpha^2 \quad (5-7)$$

Dans ces conditions pour tout ϵ positif infiniment petit, il

existe une transformation T_0 pour laquelle on a :

$$0 \leq \frac{1}{N} \sum_{i=1}^N (T_0^{-1} K' T_0^{-T})_{ii} (T^T W' T)_{ii} - M\alpha^2 < \epsilon \quad (5-8)$$

Pour trouver T_0 , il suffit de diagonaliser les matrices K et W , et d'appliquer une succession de rotations à la matrice diagonale W' afin de rendre ses éléments identiques. Le nombre de ces rotations dépend du choix de ϵ .

Le gain de bruit minimal obtenu est :

$$G_{min} = \frac{1}{N} \left[\sum_{i=1}^N \mu_i \right]^2 \quad (5-9)$$

B - Méthode de Hwang [17]

Dans cette méthode il s'agit de construire la transformation T en maintenant la contrainte de normalisation.

Soit T la transformation qui minimise le gain de bruit et qu'on factorise en :

$$T = R S \quad (5-10)$$

avec R : une matrice orthogonale ($R R^T = I$)

et S : une matrice définie positive

donc en diagonalisant S , T s'écrit :

$$T = R R_0 \Lambda R_0^T = R_1 \Lambda R_0^T \quad (5-11)$$

où $\Lambda = \text{diag} (\lambda_1 \dots \lambda_N)$

et R_0 : matrice des vecteurs propres de S , elle est orthogonale

Le gain de bruit est alors :

$$\begin{aligned} G_1 &= \text{Tr} (T^T W_0 T) = \text{Tr} (T^T T W_0) \\ &= \text{Tr} (\Lambda^2 R_1^T W_0 R_1) \\ &= \sum_{i=1}^N \lambda_i^2 r_i^2 \end{aligned} \quad (5-12)$$

où $r_i^2 = (R_1^T W_0 R_1)_{ii}$

Et T doit assurer la normalisation donc :

$$T^{-1} K_0 T^{-T} = \begin{bmatrix} 1 & X \\ X & -1 \end{bmatrix} \quad (\delta = 1) \quad (5-13)$$

Sachant que [15]

$$\det W_0 \leq \prod_{i=1}^N r_i^2 \quad (5-14)$$

et

$$\det K_0 \leq \prod_{i=1}^N \lambda_i^2$$

Et que la moyenne arithmétique est supérieure ou égale à la moyenne géométrique, ce qui permet d'écrire:

$$\frac{1}{N} \left[\sum_{i=1}^N \lambda_i^2 r_i^2 \right] \geq \left[\prod_{i=1}^N \lambda_i^2 r_i^2 \right]^{1/N} \quad (5-15)$$

alors

$$G_1 = \text{Tr}(T^T W_0 T) \geq N \left(\prod_{i=1}^N r_i^2 \right) \left(\prod_{i=1}^N \lambda_i^2 \right) \geq N |\det K_0 W_0|^{1/N} \quad (5-16)$$

Comme les valeurs propres du produit $K_0 W_0$ sont invariantes à une transformation d'état T , cette borne inférieure de (5-16) est aussi invariante.

L'égalité dans (5-16) est réalisée si et seulement si:

- 1- $K_0 W_0$ est symétrique.
- 2- K_0^{-1} et W_0 sont équivalents à une constante près.

Pour minimiser le gain :

- 1- On normalise d'abord le système par T_0 ;

$$K_1 = T_0^{-1} K_0 T_0^{-T} = I \quad \text{donc} \quad T_0 T_0^T = K_0 \quad (5-16)$$

$$W_1 = T_0^T W_0 T_0$$

- 2- On applique la matrice T décrite précédemment, donc:

$$\begin{aligned} G_1 &= \text{Tr}(T^T W_1 T) = \text{Tr}(\Lambda^2 R_1^T W_1 R_1) \\ &= \sum_{i=1}^N \lambda_i^2 r_i^2 \end{aligned} \quad (5-17)$$

Avec la contrainte

$$R_0 \Lambda^{-2} R_0^T = \begin{bmatrix} 1 & X \\ X & -1 \end{bmatrix} \quad (5-18)$$

Dans ces conditions, on obtient:

$$\sum_{i=1}^N 1/\lambda_i^2 = N \quad (5-19)$$

$$\prod_{i=1}^N \lambda_i^2 \geq 1 \quad (5-20)$$

- 3- En utilisant la fonctions de Lagrange donnée par:

$$\mathcal{L}(\lambda) = \sum_{i=1}^N \lambda_i^2 r_i^2 + \alpha \left(\sum_{i=1}^N 1/\lambda_i^2 - N \right) \quad (5-21)$$

qu'on minimise par rapport à λ_i et α pour obtenir, en considérant

la contrainte (5-19):

$$\lambda_i = \left(\frac{\sum_{i=1}^N r_i}{N r_i} \right)^{1/2} \quad (5-22)$$

donc le gain minimal est:

$$G_{\min} = \text{Tr}(T^T W_1 T) = \frac{1}{N} \left(\sum_{i=1}^N r_i \right)^2 \quad (5-23)$$

Comme

$$\det(KoWo) = \det(W_1) = \prod_{i=1}^N \mu_i^2 \quad (5-24)$$

où μ_i sont les modes du second ordre de KoWo et

$$\det(R_1^T W_1 R_1) = \det(W_1) \quad (5-25)$$

Si $R_1^T W_1 R_1$ est une matrice symétrique définie positive, alors:

$$\sum_{i=1}^N r_i^2 \geq \sum_{i=1}^N \mu_i^2 \quad (5-26)$$

avec r_i et μ_i positifs pour tout i .

L'égalité est atteinte dans (5-26) si la matrice $R_1^T W_1 R_1$ est diagonale, c'est à dire que R_1 est la matrice des vecteurs propres de W_1 .

$$\text{Donc on a : } G_{\min} = \frac{1}{N} \left(\sum_{i=1}^N \mu_i \right)^2 \geq N [\det KoWo]^{1/N} \quad (5-27)$$

La borne inférieure est atteinte si et seulement si:

$$\mu_i^2 = \mu_j^2 \quad \text{pour tous } i, j \quad (5-28)$$

4- On détermine R_0 sachant la contrainte (5-18); étant donné que R_0 est une matrice orthogonale, elle est décomposable en facteurs de matrice de rotation élémentaires de la forme:

$$R_i = \begin{bmatrix} I & \cdot & 0 & \cdot & 0 \\ \cdot & \cos\psi_i & \cdot & \sin\psi_i & \cdot \\ 0 & \cdot & I & \cdot & 0 \\ \cdot & -\sin\psi_i & \cdot & \cos\psi_i & \cdot \\ 0 & \cdot & 0 & \cdot & I \end{bmatrix} \quad i = 1, 2, \dots, N \quad (5-29)$$

(N-1) transformations sont nécessaires pour construire R_0 , qui vérifie (5-18).

5- Finalement on construit la transformation de minimisation T .

$$T = T_0 R_1 \wedge R_0^T = T_0 R_1 \wedge (R_0 \dots R_3 R_2)^T \quad (5-30)$$

qu'on applique aux paramètres du filtre (A₀, B₀, C₀, D₀).

Cette procédure est explicitée dans la section (6-2) pour le cas de filtre d'ordre 2.

5.2 Algorithmes de calcul des matrices K et W

Les matrices K et W peuvent être calculées à partir de leurs expressions:

$$K = \sum_{k=0}^{+\infty} (A^k B)(A^k B)^T = A K A^T + B B^T \quad (5-31)$$

$$W = \sum_{k=0}^{+\infty} (C A^k)^T (C A^k) = A^T W A + C^T C \quad (5-32)$$

La procédure qui calcule K calcule aussi W en remplaçant A par A^T et B par C^T. Soient les deux algorithmes suivants [1],[16]:

1^{er} Algorithme

C'est un algorithme simple et efficace, où à la i^{ème} itération, sont calculés 2ⁱ termes de la somme infinie (5-31).

Cependant, la convergence de la procédure dépend de la position des pôles du filtre par rapport au cercle unité.

La procédure est la suivante:

- (1)- Initialiser : $F \leftarrow A$ et $K \leftarrow B B^T$
- (2)- Calculer : $K \leftarrow F K F^T + K$
- (3)- Elever F au carré : $F \leftarrow F^2$
- (4)- Refaire(2) et (3) jusqu'à $F = 0$

2^{ème} Algorithme

Cet algorithme cherche d'abord une forme directe (canonique) intermédiaire, puis calcule la matrice de covariance R de cette structure qui est symétrique, et du type de Toeplitz dont la 1^{ère} ligne est [r₀, r₁, ..., r_{N-1}]. Ensuite, il calcule la matrice T qui permettra de revenir à la structure initiale en donnant la matrice K.

Les étapes de cette procédure sont les suivantes:

- (1) Pour (A,B) donnés, calculer les coefficients a_k du polynôme caractéristique : $\det(\lambda I - A) = 0$
- (2) Initialiser X(0)=0 et calculer X(k+1)= AX(k)+Ba_k pour 0 < k < N
- (3) Construire T = [X(N)...X(1)]
- (4) Calculer r(0), ..., r(N) satisfaisant l'équation suivante :

$$r(i) + \sum_{j=1}^N a_j r(|i-j|) = \begin{cases} 1 & , i=0 \\ 0 & , 0 \leq i < N \end{cases}$$

(5) Construire la matrice R (NxN) ayant la forme de Toeplitz [5]:

$$R_{ij} = r(|i-j|)$$

(6) Calculer K à partir de T et R par $K = T R T^T$

Remarque:

Le deuxième algorithme n'est pas itératif, donc il calcule K en un nombre d'étapes fini et connu à priori.

5.3 Propriétés des structures à gain de bruit minimal

5.3.1 Invariance par rapport à une transformation fréquentielle

Une propriété remarquable des structures optimales est la conservation des modes du second ordre par une transformation fréquentielle. En effet, comme les modes du second ordre d'un filtre H(z) définissent le gain en bruit minimal d'après l'équation (5-9), la famille de filtres G(z) telle que:

$$G(z) = H(F(z)) \quad (5-33)$$

Où F(z) est une transformation fréquentielle d'ordre m de la forme:

$$F(z) = \pm \prod_{i=1}^m \left(\frac{z - \alpha_i^*}{1 - \alpha_i z} \right) ; |\alpha_i| < 1 \quad (5-34)$$

et a pour modes du second ordre m duplications des modes de H(z).

En particulier, si la transformation est du type passe bas-passe bas dont l'ordre est l'unité, H(z) et G(z) ont le même gain minimal en bruit, ce qui signifie que le bruit à la sortie des structures optimales est indépendant de la largeur de la bande passante des filtres [18] (Voir chapitre VII).

Par contre, pour les structures directes le gain de bruit dépend de la bande passante du filtre du fait qu'il s'exprime en fonction des pôles de la transformation fréquentielle par [1] :

$$\text{Gain de bruit (forme directe)} = \frac{P(\alpha)}{(1 - \alpha^2)^{2N-2}} \quad (5-35)$$

où P(α) est un polynôme en α de degrés 4(N-1).

5.3.2 Expression complète du bruit de calcul

Dans certaines structures, il existe d'autres types de variables (ou noeuds) qui n'apparaissent pas dans les équations d'état du filtre (dans les structures en treillis, par exemple) et qui introduisent des erreurs de calcul supplémentaires [18].

Donc le bruit de calcul total de la sortie dépendent des N

variables d'état et des M variables correspondant aux noeuds non étatiques (qui ne sont pas issus d'un élément de retard). L'expression devient alors:

$$\sigma_{\text{Total}}^2 = \frac{\delta^2 q^2}{12} \left\{ \sum_{i=1}^N K_{ii} W_{ii} + \sum_{j=1}^M \|f_j\|_2^2 \|g_j\|_2^2 \right\} \quad (5-36)$$

(dans le cas d'un accumulateur de longueur double).

où f_i et g_i sont les réponses impulsionnelles entrée-variable d'état et variable d'état -sortie, respectivement.

5.3.3 Gain minimal et performances des réalisations

Bien que les structures optimales dans l'espace d'état offrent des gains de bruit minimums, leurs réalisations pratiques s'avèrent complexes lorsqu'il s'agit des filtres d'ordre N assez grand.

En effet, de telles structures requièrent $(N+1)^2$ multiplications pour calculer chaque échantillon de la sortie; ce qui constitue une augmentation de N^2 multiplications par rapport aux structures canoniques.

Pour réduire le nombre des multiplications dans les réalisations à bruit minimal, on utilise de multiples procédés, comme la décomposition d'un filtre d'ordre N en sections de structures optimales d'ordre 2 (ou 1) et leurs connexions en parallèle ou en cascades, ce qui réduit le nombre des multiplications à $(4N+1)$ et présente un bon compromis entre le gain de bruit de calcul et la simplicité des calculs.

D'autres procédés ont été élaborés dans le but de trouver d'autres structures de compromis, comme les procédés de tridiagonalisation de la matrice A globale, qui donnent des structures de $(5N-1)$ multiplications, les structures modulaires qui nécessitent $(3N+3)$ multiplications [19], les structures à coefficients en puissance de deux [20], etc... [21], [22], [23], [24].

Dans cette étude on s'intéressera surtout au sectionnement de l'espace d'état en cellules du second ordre connectées en parallèles.

5.4 Minimisation de la sensibilité aux coefficients

Il est important dans les considérations pratiques de synthétiser des filtres dont la fonction de transfert présente une faible sensibilité aux variations de ses coefficients dues à la limitation de leur représentation arithmétique.

Pour fonction de transfert $H(z)$

$$H(z) = \frac{\sum_{i=0}^N b_i z^{-i}}{\sum_{i=0}^N a_i z^{-i}} \quad (5-37)$$

qui est décrite par un modèle d'état contrôlable et observable

$$H(z) = D + C (zI - A)^{-1} B$$

Le fait de représenter ces coefficients par des mots de longueurs finies va faire dévier $H(z)$ de ses performances souhaitées.

Pour évaluer ces déviations dans l'espace d'état, on exprime les sensibilités de $H(z)$ par rapport à chaque composantes des paramètres d'état individuellement par [25][1]:

$$S_{a_{ij}}(z) = \frac{\partial H(z)}{\partial a_{ij}}, \quad S_{b_j}(z) = \frac{\partial H(z)}{\partial b_j}, \quad S_{c_j}(z) = \frac{\partial H(z)}{\partial c_j} \quad (5-38)$$

et qu'on peut écrire en fonction des fonctions de transfert F et G définies dans la section (4-7):

$$\begin{aligned} S_{a_{ij}}(z) &= G_i(z) F_j(z) \\ S_{b_i}(z) &= G_i(z) \\ S_{c_i}(z) &= F_i(z) \end{aligned} \quad (5-39)$$

où $G_i(z) = [C (zI - A)^{-1}]_i$

et $F_i(z) = [(zI - A)^{-1} B]_i$

et d'après les equations (4-41) et (4-55), on définit une mesure de la sensibilité par [25]:

$$M = M_A + M_B + M_C \quad (5-40)$$

avec $M_B = \frac{1}{2\pi j} \oint G^T(z) G(z) z^{-1} dz = \text{Tr}(W) = \sum_{j=1}^N W_{jj}$

$$M_C = \frac{1}{2\pi j} \oint F^T(z) F(z) z^{-1} dz = \text{Tr}(K) = \sum_{j=1}^N K_{jj}$$

$$\begin{aligned} M_A = M_B M_C &\geq \frac{1}{2\pi j} \oint \sum_{i=1}^N \sum_{j=1}^N |S_{a_{ij}}(z)|^2 z^{-1} dz \\ &= \text{Tr}(K) \text{Tr}(W) = \sum_{i=1}^N K_{ii} \sum_{j=1}^N W_{jj} \end{aligned}$$

Si on applique une transformation de normalisation telle que:

$$\text{Tr}(K') = N \quad (5-41)$$

alors l'expression de la mesure de sensibilité est:

$$\begin{aligned} M &= (N+1) \text{Tr}(W') + N \\ M &= (N+1) G' + N \end{aligned} \quad (5-42)$$

où G' est le gain de bruit de la nouvelle structure du filtre.

On remarque que la minimisation du gain entraîne celle de M donc pour un gain à minimal on a (d'après l'équation (5-9)).

$$M = \frac{(N+1)}{N} \left(\sum_{i=1}^N \mu_i \right)^2 + N \quad (5-43)$$

En conclusion, la minimisation du gain de bruit par une transformation T pour un filtre normalisé implique directement la minimisation de la mesure de la sensibilité. En conséquence, une fonction de transfert calculée à partir des coefficients d'une structure optimale normalisée présentera le moins de distorsions possibles dans le cas de l'utilisation de registres de longueurs finies pour représenter les coefficients. Cela a été confirmé pour l'exemple de simulation dans le chapitre VII.

CHAPITRE VI

STRUCTURES DECOMPOSEES OPTIMALES

STRUCTURES DECOMPOSEES OPTIMALES

6-1 Introduction

L'inconvénient majeur des structures à gain de bruit minimal étudiées dans le chapitre précédent, est la complexité de leurs réalisations matérielles due au grand nombre des multiplieurs qu'elles nécessitent. Pour un filtre d'ordre N , le nombre des multiplications par échantillons de la sortie est $(N+1)^2$ pour les structures optimales, par contre il est de $(2N+1)$ seulement pour les structures directes (canonique).

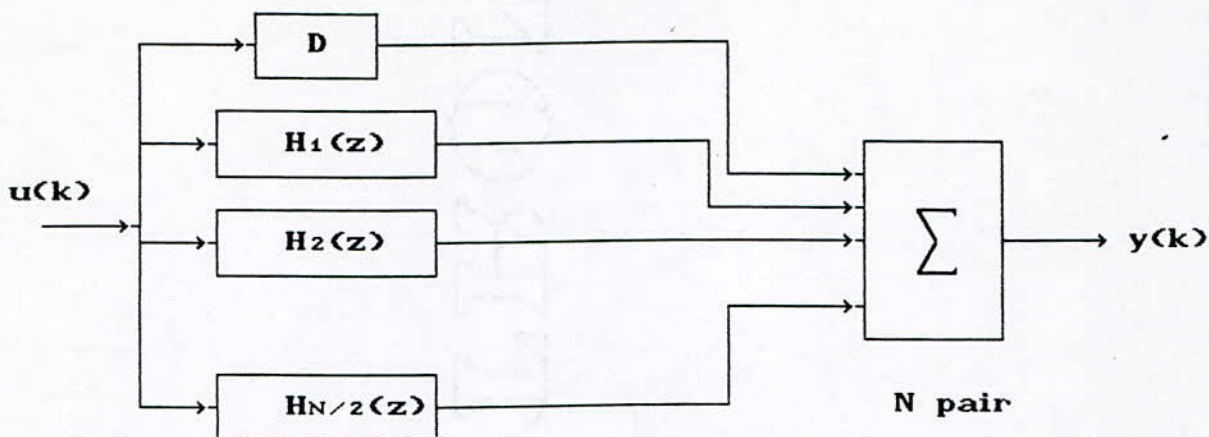
Un moyen efficace pour obtenir des structures de compromis entre le nombre des multiplications et le gain de bruit de calcul est la décomposition du filtre en cellules du second ordre regroupant chacune un couple de pôles complexes conjugués et éventuellement une cellule du premier ordre si N est impair. La connexion entre ces différentes cellules se fait soit en parallèle soit en cascade. Dans ce projet on se limite à l'étude du cas parallèle pour lequel la sortie globale du filtre est simplement la somme des sorties des différentes cellules, il en est de même pour le gain de bruit global. Dans ce cas la fonction de transfert du filtre s'écrit: (Fig(6-1))

$$H(z) = D + \sum_{i=1}^{N/2} H_i(z) \quad (\text{pour } N \text{ pair}) \quad (6-1)$$

où $H_i(z)$ est la fonction de transfert de la cellule i du second ordre de la forme:

$$H_i(z) = C_i(zI - A_i)^{-1} B_i \quad (6-2)$$

avec A_i, B_i et C_i les paramètres d'état 2×2 , 2×1 et 1×2 (resp.) et D est le chemin direct entrée-sortie du filtre .



Fig(6-1) Décomposition en cellule du second ordre d'un filtre RII

Dans la simulation des structures de compromis (chapitre VII), deux types de structures décomposées sont étudiés:

1- Structure décomposée canonique où chaque cellule a une structure canonique (ou directe) du second ordre; elle nécessite environ $(2N+1)$ multiplications par échantillon de sortie.

2- Structure décomposée optimale où chaque cellule a une structure à gain minimal, elle nécessite environ $(4N+1)$ multiplications. Bien que cette structure augmente le gain minimal du filtre, elle simplifie considérablement sa conception.

Dans ce présent chapitre, sont développés les procédés d'optimisation des cellules du second ordre.

Remarque

Il est facile de vérifier que pour la cellule du premier ordre, si elle existe (N impair), elle est optimale.

6-2 Structures optimales du second ordre

On considère le filtre d'ordre 2 dont la fonction de transfert est :

$$H(z) = C (zI - A)^{-1} B \quad (6-3)$$

où

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \quad B = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \quad C = [c_1, c_2]$$

et d'après les deux conditions nécessaires et suffisantes pour avoir une structure (A, B, C) dont la variance de l'erreur de sortie est minimale (Voir la section (5-1)).

$$C1 : K = D_0^{-1} W D_0^{-1} \quad (D_0 \text{ diagonale})$$

et

$$C2 : K_{ii} W_{ii} = K_{jj} W_{jj} \quad \text{pour tous } i,j \quad (6-4)$$

et avec la contrainte de normalisation (4-43), on a alors pour C2

$$W_{ii} = W_{jj} \quad \text{pour tous } i,j \quad (6-5)$$

et avec la condition C1, Do doit être de la forme [26]:

$$D_o = \rho I$$

$$\text{d'où} \quad W = \rho^2 K \quad (6-6)$$

Si on pose la matrice M telle que

$$M = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad (6-7)$$

on peut écrire (6-6) [26] :

$$W = \rho^2 M K M \quad (6-8)$$

En exprimant K et W en fonction des paramètres d'état (A,B,C), l'équation (6-8) entraîne:

$$A^T = M A M$$

$$C^T = \rho M B \quad (6-9)$$

En termes d'éléments, les conditions d'optimalité (6-3) deviennent [1]:

$$1- \quad a_{11} = a_{22}$$

$$2- \quad b_1 c_1 = b_2 c_2 \quad (6-10)$$

avec la contrainte de normalisation (4-43).

Il existe plusieurs structures du second ordre qui vérifient ces trois contraintes d'optimalité, parmi lesquelles apparaissent:

1- Structures de BOMAR [21]

La fonction de transfert d'un filtre du second ordre donnée par (6-2) peut aussi s'écrire

$$H(z) = \frac{q_1 z^{-1} + q_2 z^{-2}}{1 + p_1 z^{-1} + p_2 z^{-2}} \quad (6-11)$$

où q_1, q_2, p_1 et p_2 sont des constantes réelles.

Pour trouver les éléments des paramètres d'état (A, B, C) de la structure optimale correspondant à H(z); on identifie (6-11) à l'équation (6-2) développée, on obtient 4 équations à 8 inconnues:

$$\begin{aligned}
q_1 &= c_1 b_1 + c_2 b_2 \\
q_2 &= c_1 b_2 a_{12} + c_2 b_1 a_{21} - c_1 b_1 a_{22} - c_2 b_2 a_{11} \quad (6-12) \\
p_1 &= -(a_{11} + a_{22}) \\
p_2 &= a_{11} a_{22} - a_{21} a_{12}
\end{aligned}$$

Ces équations permettent d'avoir 4 réalisations équivalentes qui diffèrent seulement par les signes des coefficients $(a_{11}, a_{12}, a_{21}, a_{22}, b_1, b_2, c_1, c_2)$. Deux autres équations sont fournies par les conditions d'optimalité (6-10).

La contrainte de normalisation (4-46) pour laquelle :

$$K_{11} = K_{22} = 1/\delta^2 \quad (6-13)$$

donne à partir de la relation:

$$K = A K A^T + B B^T \quad (6-14)$$

deux autres équations.

Le système de 8 équations à 8 inconnues ainsi formé donne alors les éléments des paramètres de la structure optimale du second ordre en fonction des coefficients δ, q_1, q_2, p_1 et p_2 . La solution du système de la procédure de BOMAR est présentée par l'algorithme de calcul suivant: [1]

1) Calcul des quantités intermédiaires v_1, v_2, \dots, v_8 telles que :

$$\begin{aligned}
v_1 &= \left[\frac{q_2}{q_1} \right], \quad v_2 = (v_1^2 - p_1 v_1 + p_2), \quad v_3 = v_1 - v_2, \quad v_4 = v_1 + v_2 \\
v_5 &= p_2 - 1, \quad v_6 = p_2 + 1, \quad v_7 = v_5 \left[v_6^2 - p_1^2 \right], \quad v_8 = \left[\frac{p_1}{2} \right]^2 - p_2
\end{aligned}$$

2) Calcul des éléments des paramètres d'état

$$\begin{aligned}
a_{11} &= a_{22} = -p_1/2 \\
b_1 &= [v_7 \delta^{-2} / (2p_1 v_3 - v_6(1+v_3^2))]^{1/2} \\
b_2 &= [v_7 \delta^{-2} / (2p_1 v_4 - v_6(1+v_4^2))]^{1/2} \\
a_{21} &= [(\delta^2 b_2^2 + v_5) v_8 / (\delta^2 b_1^2 + v_5)]^{1/2} \\
a_{12} &= v_8 / a_{21} \\
c_1 &= q_1 / 2b_1 \\
c_2 &= q_1 / 2b_2
\end{aligned}$$

2- Structures de BARNES [27]

A partir de l'expression de la fonction de transfert de (6-11), on construit une structure directe (A, B, C) telle que :

4) Déterminer la matrice R_0 par la relation

$$R_0 \Lambda^{-2} R_0^T = \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}, \text{ car } K_1 = I$$

on trouve $R_0 = \begin{bmatrix} \sqrt{1/2} & -\sqrt{1/2} \\ \sqrt{1/2} & \sqrt{1/2} \end{bmatrix}$, par exemple.

En considérant le paramètre de normalisation δ et le produit ΛR_0^T

$$\text{on a : } T_1 = \delta \Lambda R_0^T = \frac{\delta}{2} \begin{bmatrix} (1 + \mu)^{1/2} & -(1 + \mu)^{1/2} \\ (1 + 1/\mu)^{1/2} & (1 + 1/\mu)^{1/2} \end{bmatrix} \text{ où } \mu = \frac{\mu_2}{\mu_1}$$

5) Appliquer T_1 à K_2 et W_2 pour obtenir K_{\min} et W_{\min} :

$$K_{\min} = \frac{1}{\delta^2} \begin{bmatrix} 1 & \frac{\mu_1 - \mu_2}{\mu_1 + \mu_2} \\ \frac{\mu_1 - \mu_2}{\mu_1 + \mu_2} & 1 \end{bmatrix}$$

$$W_{\min} = \frac{\delta^2}{4} \begin{bmatrix} (\mu_1 + \mu_2)^2 & \mu_2^2 - \mu_1^2 \\ \mu_2^2 - \mu_1^2 & (\mu_1 + \mu_2)^2 \end{bmatrix}$$

6) Calculer la transformation globale de minimisation:

$$T = T_0 R_1 T_1$$

7) Appliquer T à (A, B, C) pour avoir la structure optimale $(A_{\min}, B_{\min}, C_{\min})$. Le gain minimale est donc :

$$G_{\min} = \frac{1}{2} (\mu_1 + \mu_2)^2$$

6.3 Conclusion

Les trois procédures décrites, ci-dessus, permettent chacune un certain nombre de structures optimales du second ordre (par exemple, 4 structures pour la procédure de BOMAR) et aboutissent toutes à un même gain optimal, d'où à une même mesure de la sensibilité de la fonction de transfert du filtre globale.

Le gain total issu de toutes les cellules de structures optimales du second ordre en parallèles d'un filtre d'ordre N est:

$$G_{\text{total}} = \sum_{i=1}^{N/2} [G_{\text{min}}]_i$$

Remarque

Pour une cellule du 1^{er} ordre dont la fonction de transfert

est
$$H_1(z) = \frac{\alpha}{z-\beta}$$

le gain
$$G_1 = \frac{\alpha^2}{(1-\beta^2)^2}, \quad (\text{pour } \delta=1). \text{ Il est invariant.}$$

CHAPITRE VII

RESULTATS ET INTERPRETATIONS

RESULTATS ET INTERPRETATIONS

7.1 Introduction

Dans ce chapitre, est développée une étude comparative entre les performances de la structure canonique et de la structure décomposée optimale en parallèle d'un filtre numérique RII au moyen des simulations sur ordinateur avec l'arithmétique de la virgule fixe. Pour chaque structure, sont présentés les effets de la limitation de la précision de la représentation des coefficients du filtre et des résultats des opérations arithmétiques sur la qualité du filtrage, ainsi que le tracé de la variance de l'erreur de sortie en fonction de la longueur des mots binaires utilisée et du type d'accumulation des produits (simple ou double). De même, les tracés de cette variance sont comparés à ceux obtenus théoriquement. La sensibilité de la fonction de transfert (donc de la réponse fréquentielle) à la précision des coefficients du filtre est mise en relief par des graphes illustratifs. La variation du gain de bruit de calcul en fonction de la largeur de la bande passante du filtre est présentée afin de mettre en évidence l'effet de la transformation fréquentielle d'ordre 1 sur la performance de chaque type de structure.

Ce chapitre se termine par une comparaison globale suivie d'une conclusion.

Le filtre numérique RII simulé, dans tout le chapitre, est un filtre passe-bas de BUTTERWORTH d'ordre 8 dont la fréquence de coupure normalisée est $f_c = 0.125$ ($\omega_c = 0.25 \pi$).

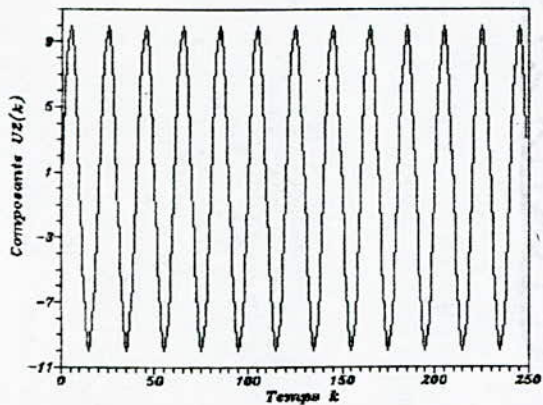
Les signaux d'entrée générés et utilisés sont:

- Un signal S_1 qui est la somme de deux sinusoides dont l'une des fréquences se trouve dans la bande passante du filtre et l'autre dans sa bande d'affaiblissement (voir Fig-(7-1-a,b,c)).

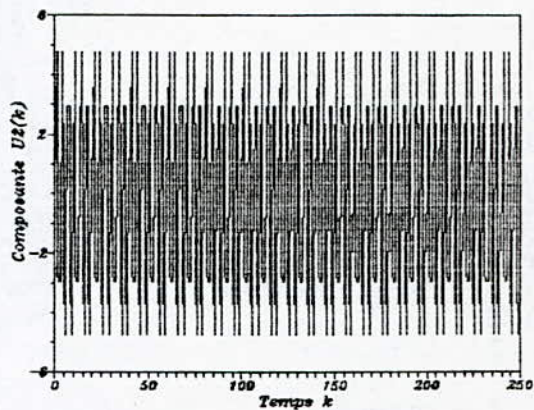
Ce signal est donné par:

$$U(k) = 10 \sin(2 \pi f_1 k) + 5 \sin(2 \pi f_2 k) \quad (7- 1)$$

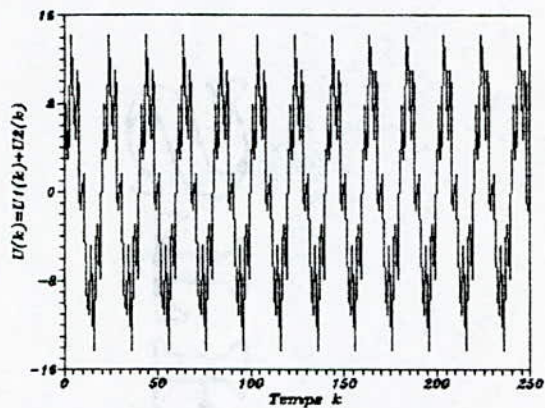
avec $f_1 = 0.05$ ($\omega_1 = 0,1 \pi$) et $f_2 = 0.3$ ($\omega_2 = 0.6 \pi$)



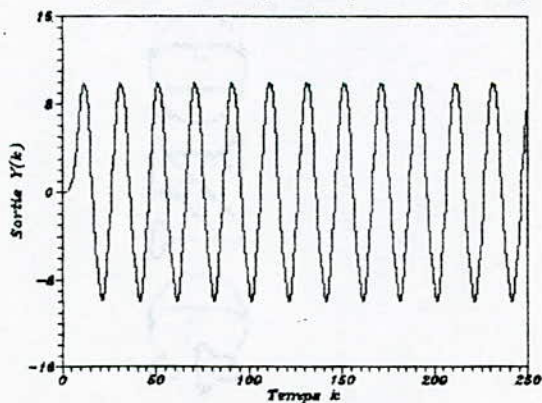
(a) COMPOSANTE DE FREQUENCE NORMALISEE 0.05



(b) COMPOSANTE DE FREQUENCE NORMALISEE 0.3

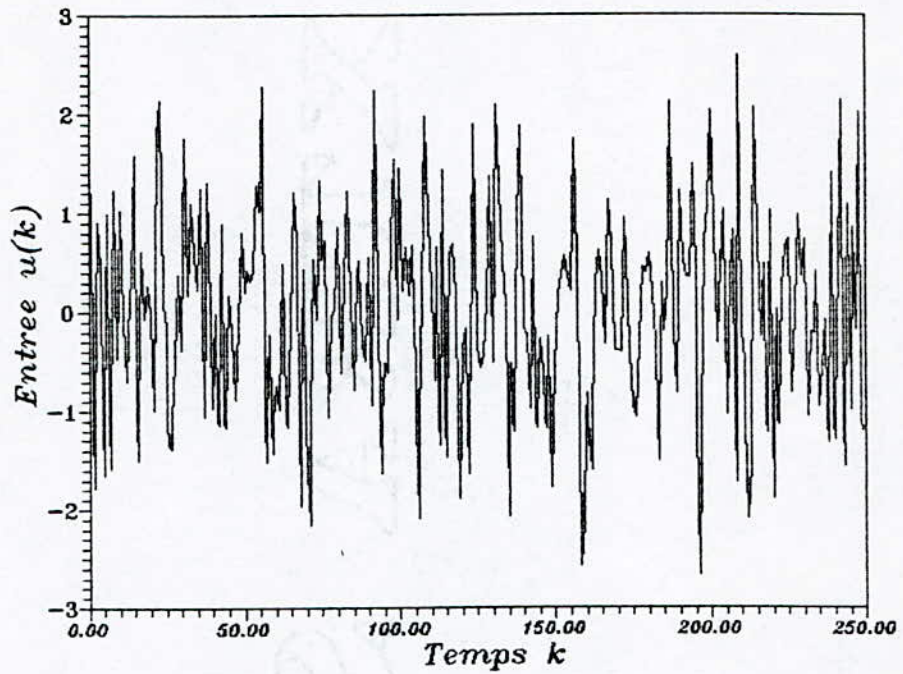


(c) Signal d'entree du filtre
(somme des composantes $f_1=0.05$ et $f_2=0.3$)

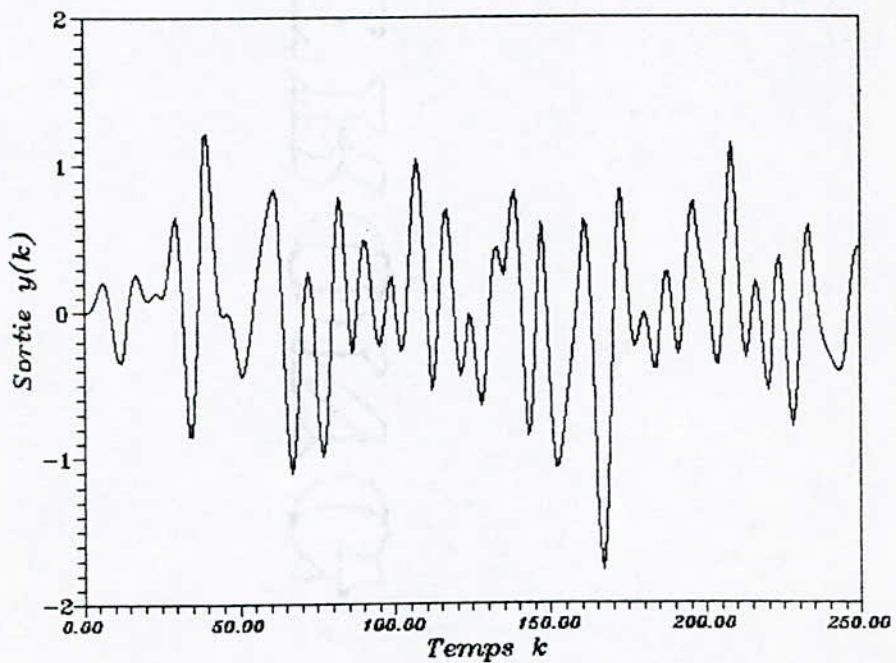


(d) signal de sortie (54 bits)

FIG-(7-4) FILTRAGE D'UN SIGNAL SOMME DE 2 SINUSOIDES
AVEC UNE PRECISION INFINIE (54 Bits)



(a) Signal non filtré



(b) Signal Filtré

FIG-(7-2) FILTRAGE D'UN SIGNAL BRUIT BLANC
CENTRE DE VARIANCE UNITE

- Un signal S2 aléatoire de distribution gaussienne centrée de variance l'unité simulant un bruit blanc (voir Fig(7-2-a)).

7.2 Performances d'un filtre numérique avec une structure canonique

La fonction de transfert du filtre utilise est donnée par:

$$H(z) = \frac{\sum_{i=0}^8 b_i z^{-i}}{1 + \sum_{i=1}^8 a_i z^{-i}} \quad (7-2)$$

dont les coefficients du numérateur sont:

$b_0 = 1.0791128792196987E-04$
 $b_1 = 8.6329030337575899E-04$
 $b_2 = 3.0215160618151565E-03$
 $b_3 = 6.0430321236303130E-03$
 $b_4 = 7.5537901545378912E-03$
 $b_5 = 6.0430321236303130E-03$
 $b_6 = 3.0215160618151565E-03$
 $b_7 = 8.6329030337575899E-04$
 $b_8 = 1.0791128792196987E-04$

et les coefficients du dénominateur sont:

$a_0 = 1.0000000000000000$
 $a_1 = -3.983784327663337$
 $a_2 = 7.536234286077626$
 $a_3 = -8.599815333907508$
 $a_4 = 6.400154303626428$
 $a_5 = -3.156025399574342$
 $a_6 = 1.001696629290593$
 $a_7 = -0.1863424879766865$
 $a_8 = 1.5507616199615062E-02$

7.2.1 Simulation avec précision infinie

Le filtrage des signaux S1 et S2 s'effectue suivant la relation d'entrée-sortie suivante:

$$y(k) = \sum_{i=0}^8 b_i u(k-i) - \sum_{i=1}^8 a_i y(k-i) \quad (7-3)$$

La représentation binaire des coefficients et celle des résultats des opérations (multiplications et additions) est faite sur des mots de 54 bits (précision de l'ordinateur) pour lesquels la précision est supposée infinie.

En effet, on constate que le filtrage de S1 est correctement effectué (voir Fig(7-1-d)); la fréquence f2 qui se trouve à l'extérieur de la bande passante est bien rejetée, seule la fréquence f1 est retenue à la sortie.

Le filtrage du signal bruit blanc S2 élimine les composantes de hautes fréquences c'est à dire celles qui sont supérieures à fc. La sortie du filtre pour ce signal est donnée par Fig(7-2-b).

7.2.2 Simulation avec précision finie

Le filtrage est réalisé à partir de la même relation (7-3), mais cette fois-ci en limitant la longueur des mots binaires représentant le signal de l'entrée, les coefficients du filtre et les registres de stockage et d'accumulation. Cette quantification est réalisée par la fonction BBIT dont le programme en Fortran est donné par Fig(7-3). Cette fonction effectue la quantification par l'arrondi avec saturation. La mobilité de la virgule dans cette quantification est destinée à élargir la dynamique afin de ne considérer que les erreurs de calcul, (la probabilité de dépassement étant alors faible). Pour simuler le quantificateur réel en virgule fixe, on propose une autre fonction Fortran, BIT, donnée par figure Fig(7-3 bis):

On effectue le filtrage pour des longueurs de mots fixées à B=4, 8, 12 et 16 bits. La sortie obtenue pour les quatre cas, est présentée par Fig(7-4).

On constate que le filtrage devient erroné à partir de 12 bits; la sinusoïde de la sortie est distordue et cela revient au fait que les coefficients deviennent de moins en moins représentables. Pour 4 bits le signal en sortie est nul car tous les coefficients du numérateur de la fonction de transfert sont représentés par la valeur zero.

7.2.3 Variance de l'erreur de sortie

Pour chaque valeur de B, est calculée l'erreur de la sortie par l'expression :

$$\varepsilon(k) = Y_i(k) - Y_f(k) \quad (7-4)$$

où $Y_i(k)$: la sortie du filtre avec précision infinie,

$Y_f(k)$: la sortie du filtre avec précision finie.

On calcule la variance de l'erreur de sortie simulée (qu'on

La fonction "BBIT" simule un quantificateur par l'arrondi avec une caractéristique de dépassement par saturation où:

IB= la longueur du mot désirée (bit de signe exclu).

X = la valeur à quantifier (en double précision).

XX= X quantifié sur IB+1 bits.

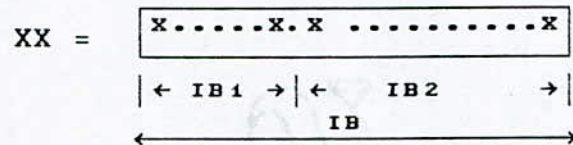
```
REAL*8 function BBIT(X,IB)
REAL*8 A,y,xx,k,k1,N,J,X
y= Dabs(x)
A= INT(y)
K= 2** IB
J= K - 1
IF ( A .GE. J ) THEN
    XX= J
    GOTO 20
ENDIF
IF ( A .EQ. 0.DO ) THEN
    xx= NINT( y * K ) / K
    GOTO 20
ENDIF
I= -1
10 K= K / 2
I= I + 1
N= INT(A / K)
IF ( N .EQ. 0 ) GOTO 10
K1= 2 * I
XX = NINT( (y-A)* K1 ) / K1 + A
20 IF (x.LT.0.d0) xx = -xx
bbit= xx
RETURN
END
```

FIG-(7- 3) FONCTION DE QUANTIFICATION " BBIT "

La fonction BIT réalise la quantification en virgule fixe par l'arrondi en signe et valeur absolue et une caractéristique de dépassement par saturation ; où :

x est la valeur à quantifier (précision infinie) ,
 IB est la taille du mot en bits ,
 IB2 est la longueur en bits de la partie fractionnaire du mot.

XX est la valeur quantifiée de X qui sera sous forme:



```

REAL*8 FUNCTION BIT(X,IB,IB2,XX)
REAL*8 y,xx,k1,X,AX,S
INTEGER IB,IB1,IB2
IB2= IB - IB1
AX = 0.d0
DO 5 I=1,IB2
  AX= AX + 2.d0**(FLOAT(-I))
5  CONTINUE
Y=DABS(X)
K1= 2.d0**FLOAT(IB2)
S= (2.d0**FLOAT(IB1))- 1.d0+ AX
IF( Y. GE. S) THEN
  XX = S
  Write(*,*)' OVERFLOW '
  GOTO 20
ENDIF
XX= NINT(Y * k1) / k1
20 IF ( X .LT. 0.d0) XX = - XX
BIT = XX
RETURN
END
  
```

Fig(7-3 bis) Fonction de quantification en virgule fixe

appellera pratique) par:

$$\sigma_y^2 \text{pratique} = \frac{1}{N_e} \sum_{k=0}^{N_e-1} \varepsilon(k)^2 - \left[\frac{1}{N_e} \sum_{k=0}^{N_e-1} \varepsilon(k) \right]^2 \quad (7-5)$$

où N_e est le nombre d'échantillons du signal de sortie.

pour deux cas:

- Cas de l'accumulateur simple, pour lequel chaque produit intermédiaire est quantifié avant la sommation dans (7-3).

L'expression théorique de la variance des erreurs de calcul en sortie est:

$$\sigma_y^2 \text{théorique} = \frac{q^2}{12} (N+1) \delta^2 (G + 1) \quad (7-6)$$

N est l'ordre du filtre, dans ce cas $N = 8$

q est le pas de quantification, dans ce cas $q = 10 \cdot 2^{(-B+1)}$

δ est le paramètre de normalisation, ici $\delta = 1$,

G est le gain de bruit de calcul de la structure canonique, il a été calculé par un programme en Fortran (établi par B.DERRAS à l'ENP en 1991), et pour le filtre utilisé, sa valeur est:

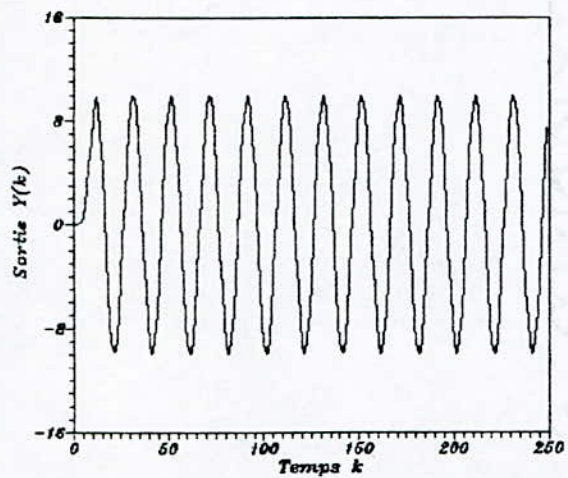
$$G_{\text{canonique}} = 8019.683673775$$

- Cas de l'accumulateur double, pour lequel la quantification s'effectue après la sommation de tous les produits de (7-3).

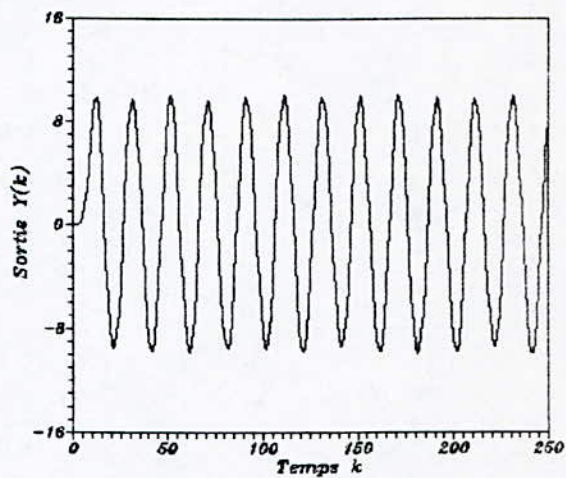
L'expression théorique de la variance, dans ce cas, est:

$$\sigma_y^2 \text{théorique double} = \frac{\sigma_y^2 \text{théorique simple}}{(N+1)} \quad (7-7)$$

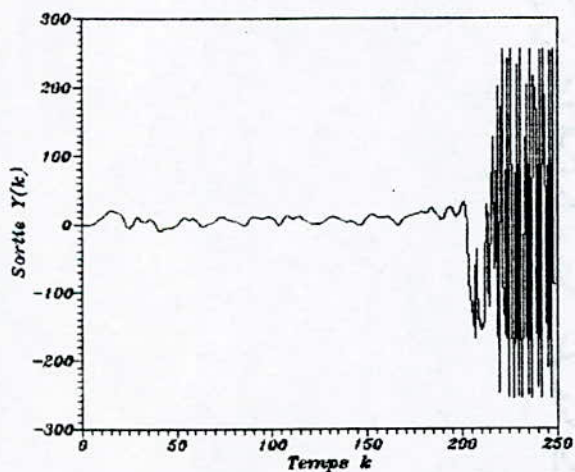
La variance de l'erreur est tracée en fonction de B pour les signaux S1 et S2 (voir Fi (7-5) et Fig(7-6)). On constate que, pour ces deux tracés et dans tous les cas, la variance est décroissante et correspond bien à aux expressions théoriques données par (7-6) et (7-7). On peut aussi remarquer que la variance pour le cas d'un accumulateur double est inférieure à celle pour un accumulateur simple. La variance relative au bruit blanc S2 se rapproche le plus de la variance théorique que celle relative au signal S1.



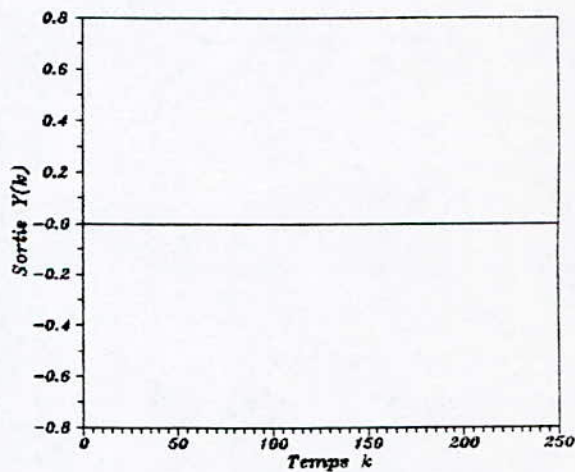
(a) Pour Accumulateur 18 BITS



(b) Pour Accumulateur 12 BITS

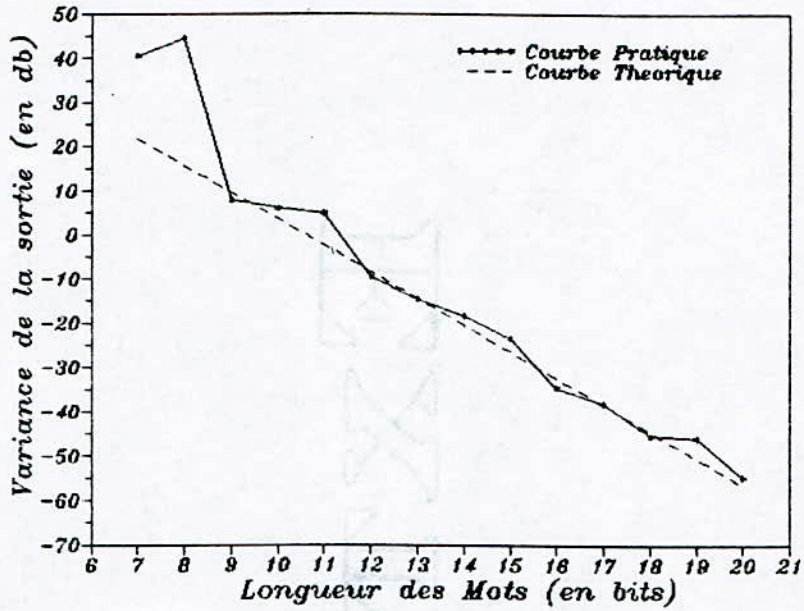


(c) Pour Accumulateur 8 BITS

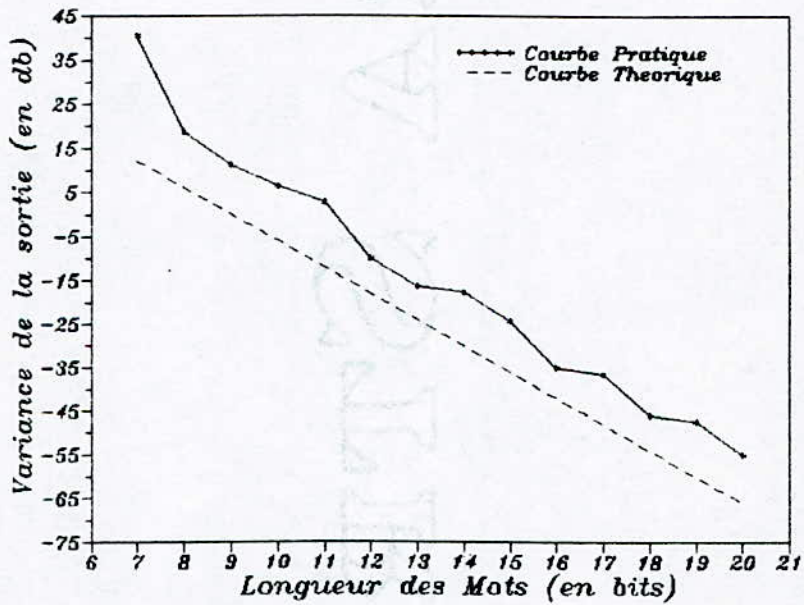


(d) Pour Accumulateur 4 BITS

FIG-(7-4) SORTIE DE LA STRUCTURE CANONIQUE
POUR DES DIFFERENTES LONGUEURS D'ACCUMULATEURS

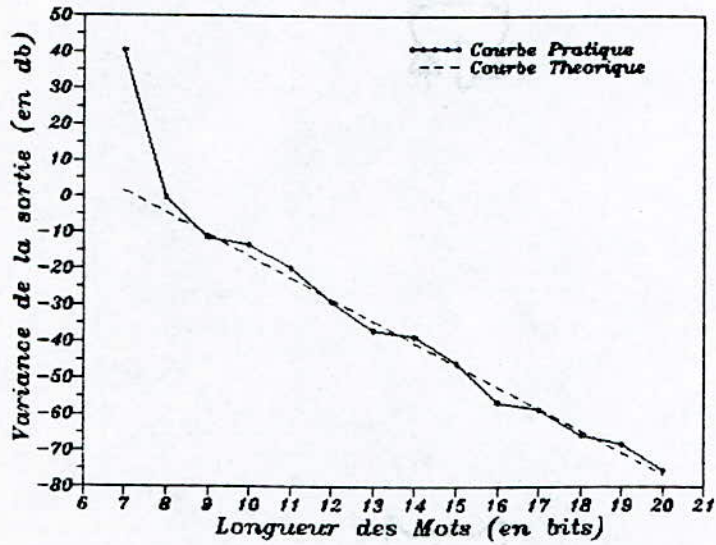


(a) Accumulateur Simple

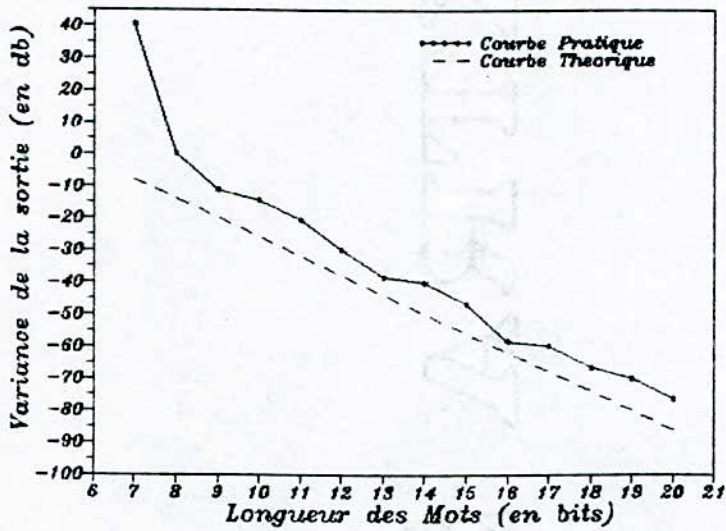


(b) Accumulateur Double

FIG-(7-5) ENTREE SOMME DE DEUX SINUSOIDES
($f_1 = 0.05$ et $f_2 = 0.3$)



(a) Accumulateur Simple



(b) Accumulateur Double

FIG-(7-6) ENTREE BRUIT BLANC CENTREE DE VARIANCE UNITE

appellera pratique) par :

$$\sigma_y^2 \text{pratique} = \frac{1}{N_e} \sum_{k=0}^{N_e-1} \varepsilon(k)^2 - \left[\frac{1}{N_e} \sum_{k=0}^{N_e-1} \varepsilon(k) \right]^2 \quad (7-5)$$

où N_e est le nombre d'échantillons du signal de sortie.

pour deux cas :

- Cas de l'accumulateur simple, pour lequel chaque produit intermédiaire est quantifié avant la sommation dans (7-3).

L'expression théorique de la variance des erreurs de calcul en sortie est :

$$\sigma_y^2 \text{théorique} = \frac{q^2}{12} (N+1) \delta^2 (G + 1) \quad (7-6)$$

N est l'ordre du filtre ,dans ce cas $N = 8$

q est le pas de quantification, dans ce cas $q = 10 \cdot 2^{-(B+1)}$

δ est le paramètre de normalisation, ici $\delta = 1$,

G est le gain de bruit de calcul de la structure canonique, il a été calculé par un programme en Fortran (établi par B.DERRAS à l'ENP en 1991), et pour le filtre utilisé, sa valeur est :

$$G_{\text{canonique}} = 8019.683673775$$

- Cas de l'accumulateur double, pour lequel la quantification s'effectue après la sommation de tous les produits de (7-3). L'expression théorique de la variance, dans ce cas, est :

$$\sigma_y^2 \text{théorique double} = \frac{\sigma_y^2 \text{théorique simple}}{(N+1)} \quad (7-7)$$

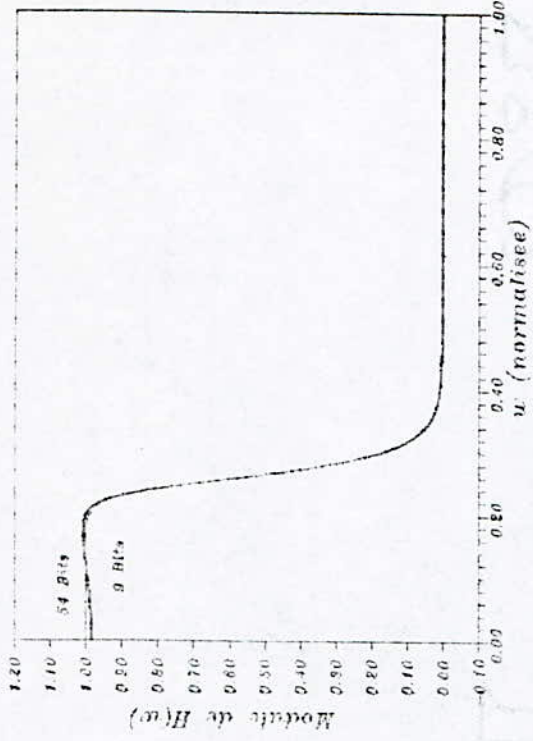
La variance de l'erreur est tracée en fonction de B pour les signaux S1 et S2 (voir Fi (7-5) et Fig(7-6)). On constate que, pour ces deux tracés et dans tous les cas, la variance est décroissante et correspond bien à aux expressions théoriques données par (7-6) et (7-7). On peut aussi remarquer que la variance pour le cas d'un accumulateur double est inférieure à celle pour un accumulateur simple. La variance relative au bruit blanc S2 se rapproche le plus de la variance théorique que celle relative au signal S1.

7.2.4 Sensibilité de la réponse fréquentielle

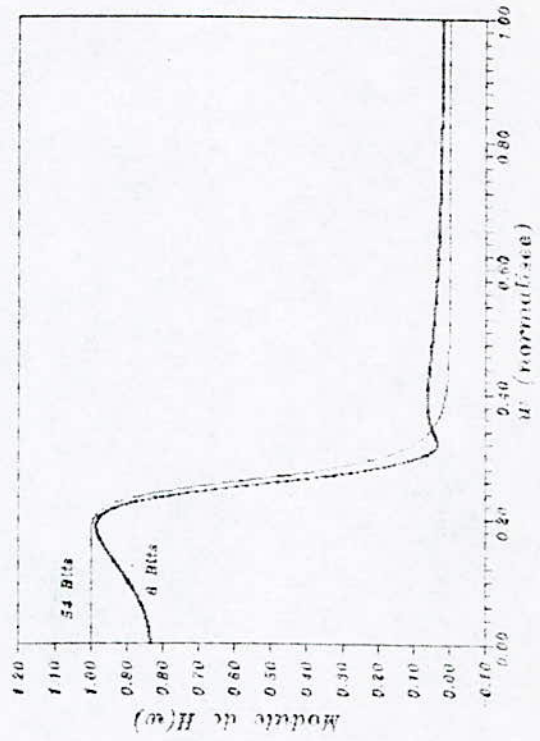
L'effet de la variation des coefficients de la structure canonique due à la limitation de la longueur des registres sur la fonction de transfert du filtre est présenté par Fig (7-7).

Cette figure montre que, pour 12 bits, une distorsion de la

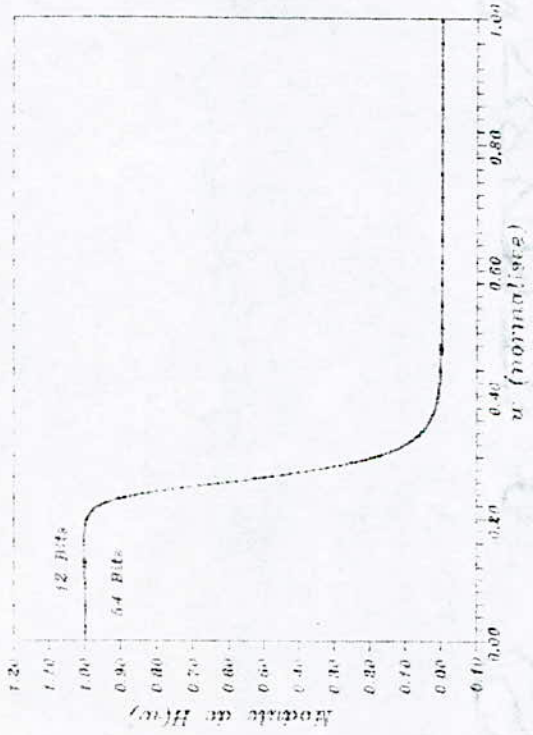
Filter de BUTTERWORTH d'ordre 8



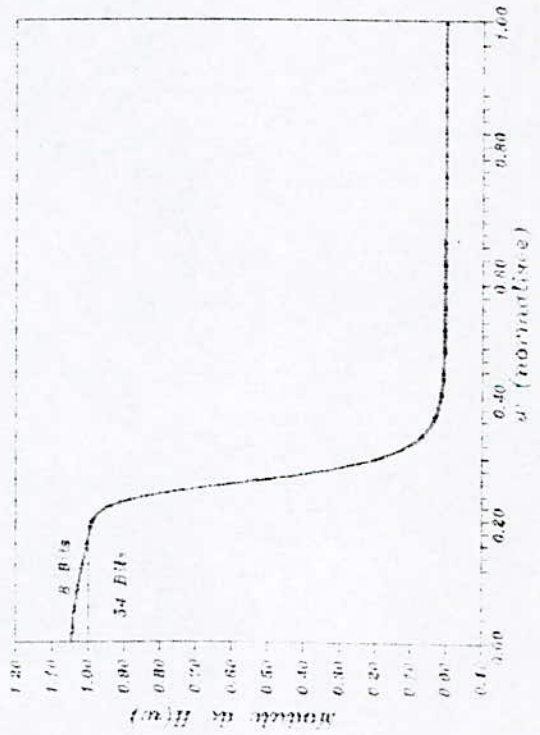
Filter de BUTTERWORTH d'ordre 8



Filter de BUTTERWORTH d'ordre 8



Filter de BUTTERWORTH d'ordre 8



Fig(7-12)

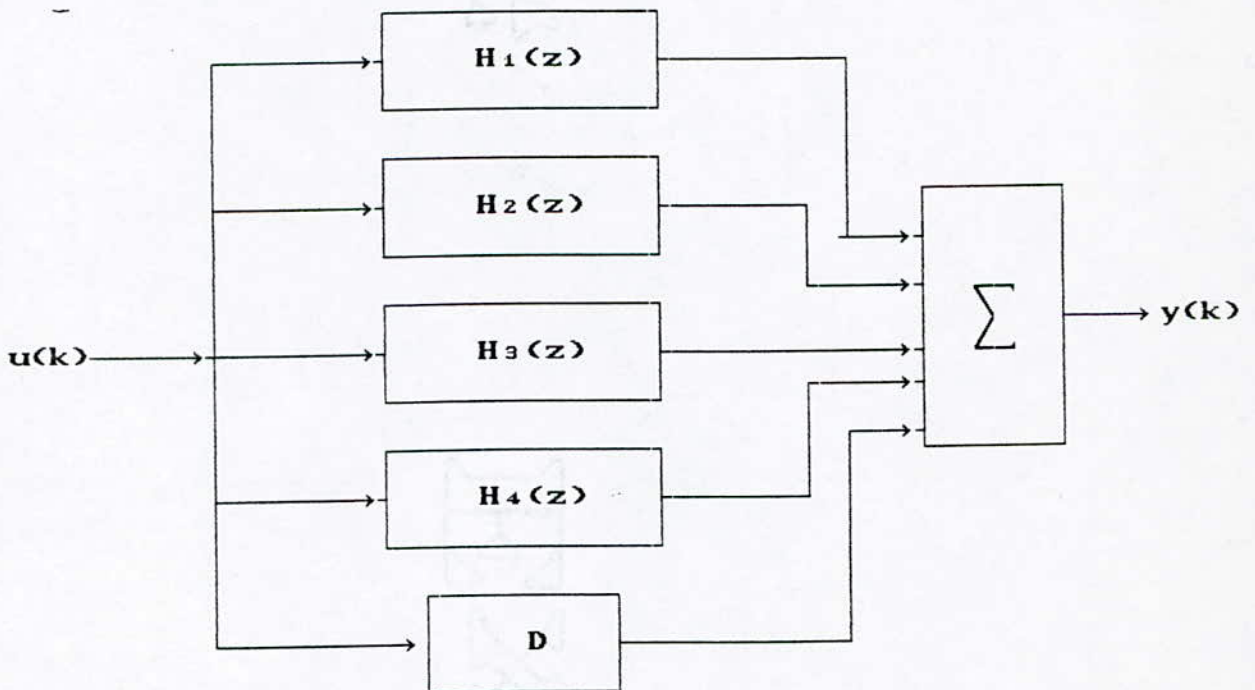
fonction de transfert apparait. Elle devient plus importante au fur et à mesure que la longueur des registres contenant les coefficients devient petite.

Pour 6 bits, la fonction de transfert est nulle car les coefficients ne sont plus représentables dans cette gamme. En conséquence, la qualité du filtrage se détériore graduellement pour des registres moins de 12 bits. La sensibilité de la fonction de transfert se traduit par l'ampleur de cette dégradation.

7.3 Performances d'un filtre numérique avec une structure décomposée en parallèle

Dans cette partie, il s'agit de mettre en œuvre la théorie explicitée dans le chapitre précédent par une simulation sur ordinateur. A cet effet, la décomposition du filtre en sections du second ordre regroupant chacune un couple de pôles complexes conjugués et disposées en parallèle, ainsi que la construction des structures canoniques et optimales suivant les trois méthodes exposées dans le chapitre V (BOMAR, BARNES et HWANG), sont exécutées par le programme STRUCMIN se trouvant en annexe.

Le filtre d'ordre 8 simulé est ainsi décomposé en 4 cellules d'ordre 2 que l'on connectera en parallèle (Fig (7-8)).



Fig(7-8) Décomposition du filtre d'ordre 8 en cellules parallèles

Chaque cellule est de la forme :

$$H(z) = \frac{q_1 z + q_2}{1 + p_1 z + p_2 z^2} \quad (7-8)$$

où q_1, q_2, p_1 et p_2 sont des constantes réelles spécifiques pour chaque cellule. Les pôles de $H(z)$ sont de la forme:

$$\lambda = \alpha \pm j \beta$$

Les résultats de la décomposition sont donnés par le tableau suivant :

	$H_1(z)$	$H_2(z)$	$H_3(z)$	$H_4(z)$
q1	-3.2868476268	-0.6181903311	0.4356772381	3.4706539054
q2	1.9459927643	-0.4018617159	-0.3054473169	-1.1160909524
p1	-0.8905976013	-1.0153398695	-1.2427733740	-0.8350734828
p2	0.2594951881	0.4359073931	0.7575469450	0.1809722346
α	0.4452987909	0.5076699257	0.6213867068	0.4175367355
β	0.2473947555	0.4221121073	0.6094468832	0.0814573765

Le coefficient du chemin direct entrée/ sortie est:

$$D = 1.0791128792196987E-04$$

On construit les structures d'états directes (canoniques) des 4 cellules et en calculant le gain de bruit pour chacune, on obtient:

Cellule	1		2		3		4	
Matrice A	0.000	1.000	0.000	1.000	0.000	1.000	0.000	1.000
	-0.259	0.891	-0.436	1.015	-0.757	1.243	-0.181	0.835
Vecteur B	0.000		0.000		0.000		0.000	
	1.000		1.000		1.000		1.000	
Vecteur C	1.946	-3.287	-0.402	-0.618	-0.305	0.436	-1.116	3.471
Gain de bruit	35.334463		6.892163		3.728628		37.074119	

Le gain total pour la forme décomposée directe est donc :

$$G_{\text{déc. can}} = 83.0293732064$$

On calcule ensuite les structures d'état optimales selon une des trois procédures étudiée dans le chapitre VI et on procède de la même manière que pour les structures canoniques. En utilisant la procédure de BOMAR, par exemple, on obtient les résultats suivant:

Cellule	1		2		3		4	
Matrice A	0.445	-0.486	0.508	-0.279	0.621	-0.700	0.417	-0.035
	0.126	0.445	0.639	0.508	0.530	0.621	0.191	0.417
Vecteur B	0.938		0.886		0.495		0.923	
	0.839		0.237		0.491		0.797	
Vecteur C	-1.752	-1.958	-0.349	-1.304	0.440	0.444	1.879	2.178
Gain de bruit	8.7138762		4.327928		1.610885		11.126727	

Le gain total de la structure décomposée optimale est :

$$G_{\text{déc. optimale}} = 25.77941741187027$$

De même , il est intéressant de voir le gain de bruit de la structure optimale globale du filtre à partir de ses modes du second ordre et dont la valeur est [19] :

$$G_{\text{optimal}} = 0.8597278167$$

D'après ces résultats, on déduit les hiérarchies qui existent entre les gains des différentes structures :

$$G_{\text{opt. global}} < G_{\text{déc. optimal}} < G_{\text{déc. can}} < G_{\text{can. global}}$$

et entre les nombres d'opérations NO par échantillon de la sortie pour chaque structure :

$$NO_{\text{opt. global}} > NO_{\text{déc. optimal}} > NO_{\text{déc. can}} = NO_{\text{can. global}}$$

où $NO_{\text{opt. global}} = (N+1)^2$, $NO_{\text{déc. optimal}} = 4N+1$ et $NO_{\text{déc. can}} = 2N+1$

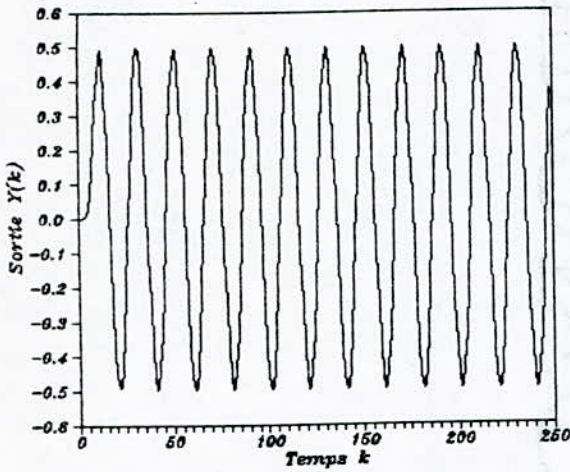
On peut conclure que les structures décomposées canonique et optimale présentent un certain équilibre par rapport aux structures optimale et canonique globales du point de vue des erreurs de calcul et nombres d'opérations. Le choix de la structure convenable dépend du cahier de charge exigé pour une application donnée, et plus précisément du matériel disponible (multiplieurs), du coût de la réalisation et du rapport signal sur bruit désiré (qualité du filtrage).

7.3.1 Simulation avec précision finie

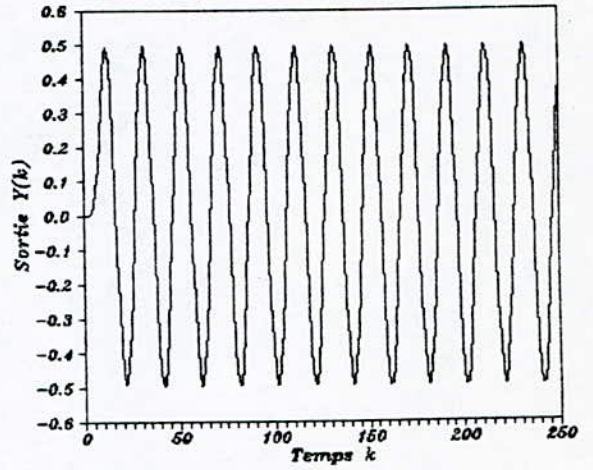
Dans cette partie ,on refait la même opération de filtrage sur le signal S1 qu'on normalise en le mettant dans la gamme [-1,1]. on fait varier la longueur des mots binaires mais cette fois avec une structure décomposée optimale; On obtient la Fig(7-9) sur laquelle on remarque que la sinusoïde de la sortie est conservée jusqu'à environ 8 bits, et que pour 6 bits le signal n'est plus nul bien qu'il est distordu relativement au cas canonique (voir Fig(7-4)).

a- variance de l'erreur de sortie

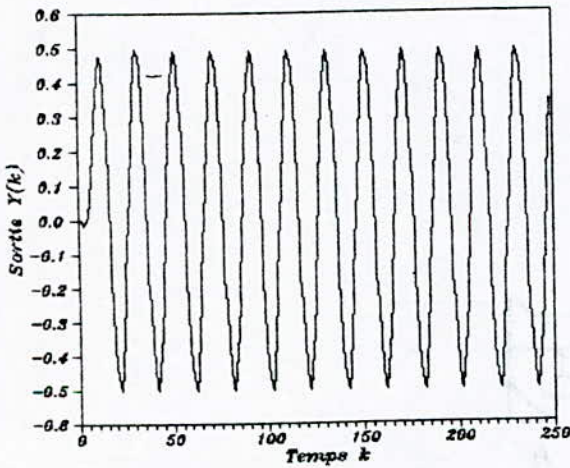
De la même manière que pour le cas de la structure canonique, on calcule la variance de l'erreur de sortie en fonction du nombre de bits des registres simulés , on obtient pour le signal S1, la figure Fig(7-11) et pour le signal S2 la figure Fig (7-10).



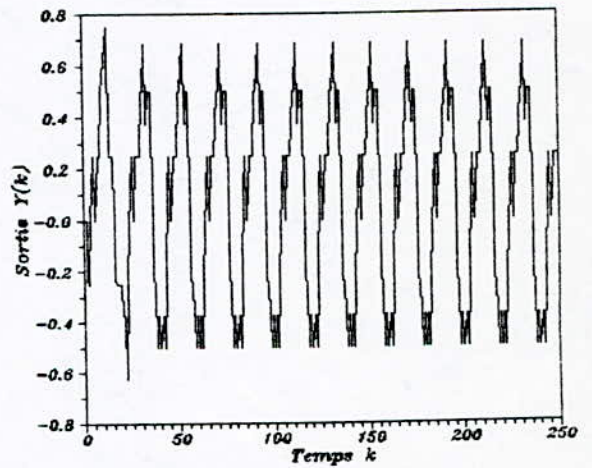
(a) Pour Accumulateur 16 BITS



(b) Pour Accumulateur 12 BITS



(c) Pour Accumulateur 8 BITS



(d) Pour Accumulateur 4 BITS

FIG-(7- 9) SORTIE DE LA STRUCTURE DECOMPOSEE OPTIMALE OPTIMALE POUR DES DIFFERENTES LONGUEURS D'ACCUMULATEURS

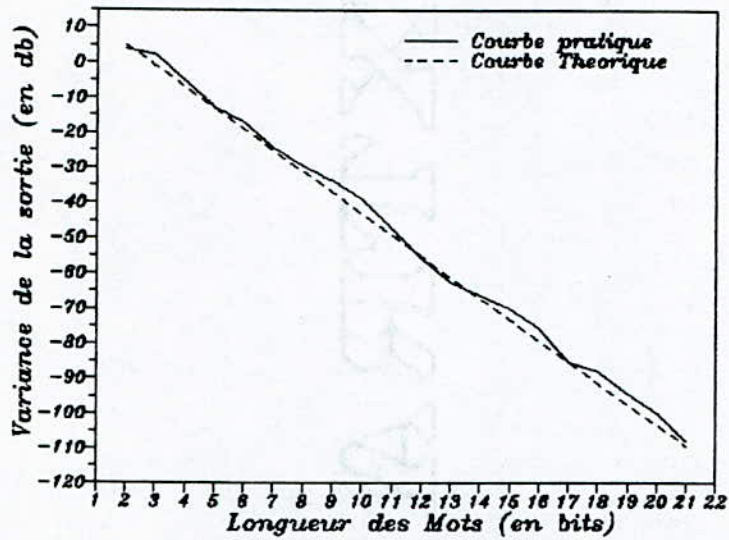


FIG-(7-10 a) ENTREE BRUIT BLANC
Accumulateur SIMPLE

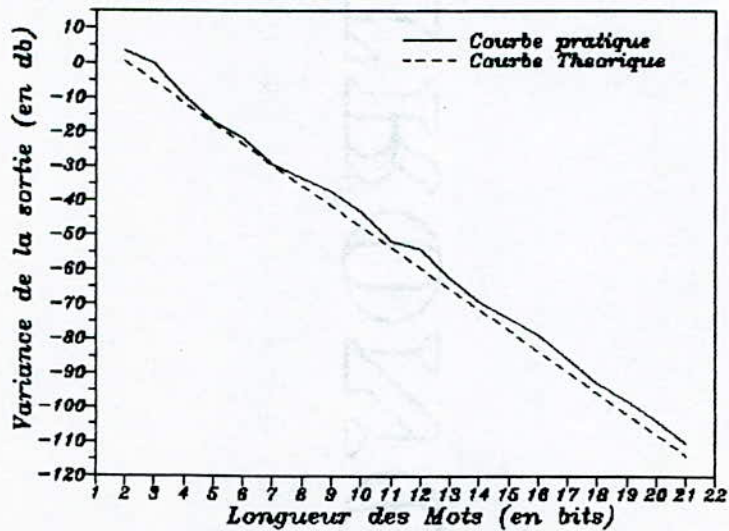


FIG-(7-10 b) ENTREE BRUIT BLANC
Accumulateur de longueur DOUBLE

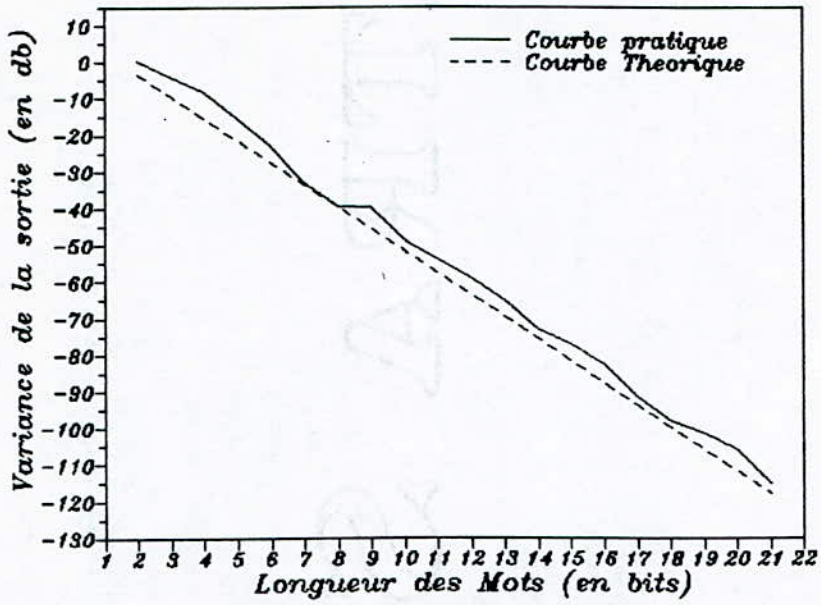


FIG-(7-11 a) ENTREE SINUSOIALE
Accumulateur SIMPLE

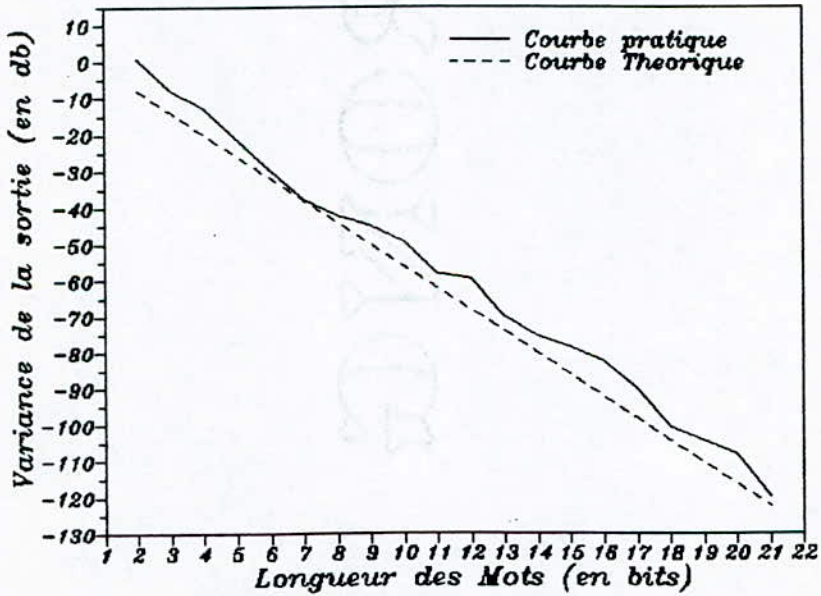


FIG-(7-11 b) ENTREE SINUSOIALE
Accumulateur de longueur DOUBLE

On remarque pour chaque signal le rapprochement entre la variance théorique donnée par :

$$\sigma_y^2 \text{ théorique} = \frac{q^2}{12} (3 \delta^2 G_{\text{déc. opt}} + 1) \quad (7-8)$$

pour le cas d'un accumulateur simple et par :

$$\sigma_y^2 \text{ théorique} = \frac{q^2}{12} (\delta^2 G_{\text{déc. opt}} + 1) \quad (7-10)$$

pour le cas d'un accumulateur double; et entre la variance pratique . Ce qui confirme l'applicabilité de la théorie de minimisation par les structures d'état de MULLIS et de HWANG[16],[17]. Pour le cas de l'accumulateur double, on remarque que la variance est moindre que celle du cas de l'accumulateur simple.

b- Sensibilité de la réponse fréquentielle

En utilisant les paramètres de la structure décomposée optimale, on calcule pour chaque cellule la réponse fréquentielle par:

$$H_i (e^{j\omega}) = C_i (zI - A_i)^{-1} B_i \quad (7-11)$$

où A_i , B_i , C_i sont les paramètres d'état de la cellule i d'ordre 2.

Ensuite on calcule la réponse fréquentielle globale :

$$H (e^{j\omega}) = \sum_{i=1}^4 H_i (e^{j\omega}) + D \quad (7-12)$$

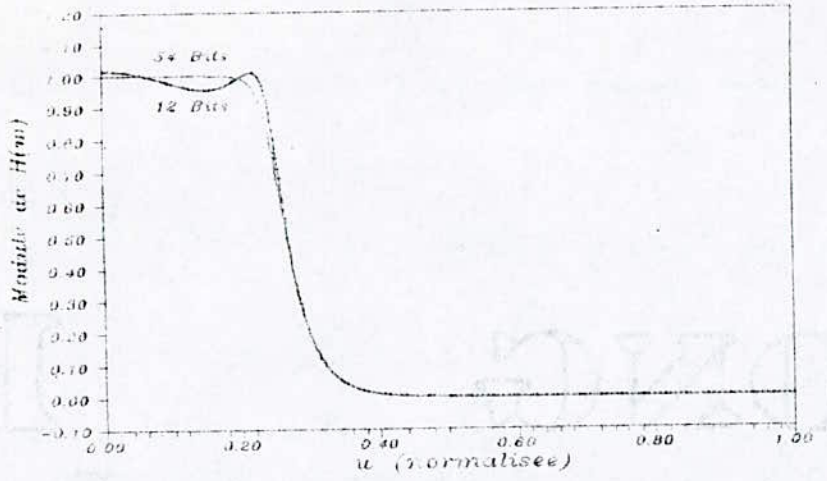
On obtient pour les différentes valeurs de longueur des registres, les courbes données par Fig(7-12).

On constate une nette amélioration de la réponse fréquentielle pour cette structure en comparaison avec la structure canonique en précision finie Fig(7-7)) . En effet , la sensibilité de la fonction de transfert a diminué; cela revient au fait que la décomposition en sections d'ordre 2 a permis de séparer les pôles qui sont proches dans le cercle unité , avant la quantification.

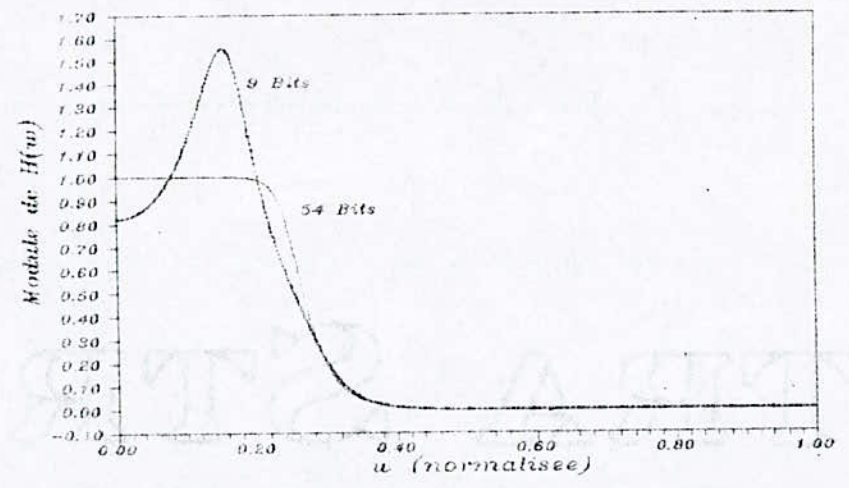
La distorsion de la courbe commence à apparaître à partir de 8 bits. Cette réponse n'est plus nulle pour 6 bits comme il a été vu pour la structure canonique globale; (comparer Fig(7-7)et Fig(7-12)).

Fig (7-7)

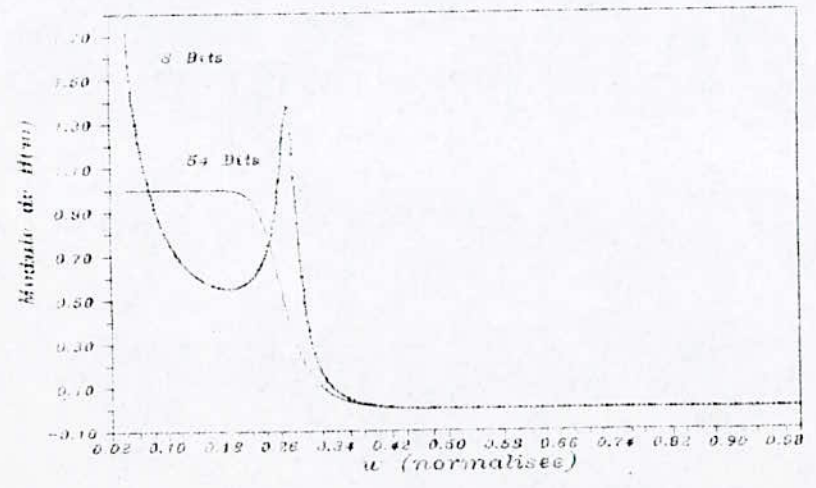
Filtre de BUTTERWORTH d'ordre 8



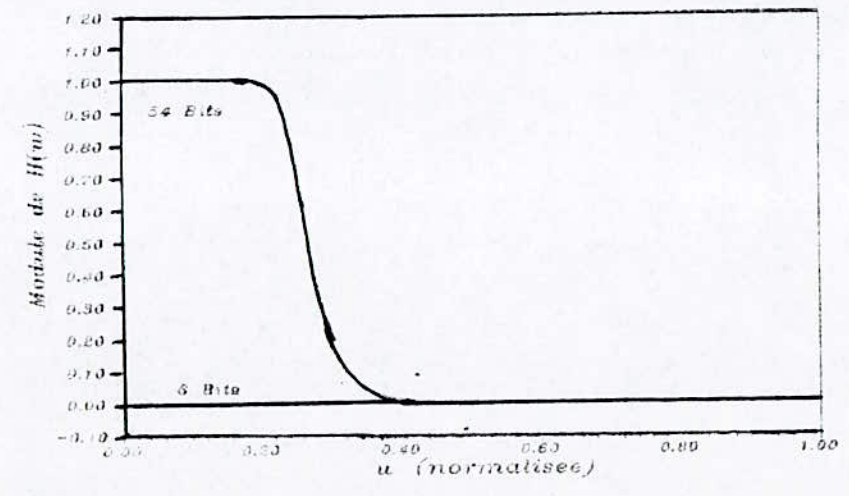
Filtre de BUTTERWORTH d'ordre 8



Filtre de BUTTERWORTH d'ordre 8



Filtre de BUTTERWORTH d'ordre 8



7.4 Effet d'une transformation fréquentielle sur le gain de bruit de calcul

Dans cette section, On étudie le comportement des gains de bruit de calcul des quatre types de structures vis-à-vis - d'une transformation fréquentielle d'ordre 1, en synthétisant des filtres passe-bas de BUTTERWORTH d'ordre 8 de fréquences de coupure différentes. On obtient le tableau suivant:

Fréquence de coupure	Gain de Bruit de Calcul pour la Structure :			
	Canonique	Décomposée Canonique	Décomposée Optimale	Optimale
0.1	3.6697E+009	434.747	25.7794	0.859728
0.2	196390	120.161	25.7794	0.859728
0.3	590.243	63.4288	25.7794	0.859728
0.4	12.7886	45.8975	25.7794	0.859728
0.5	10.92	41.5147	25.7794	0.859728
0.6	157.627	45.8975	25.7794	0.859728
0.7	8697.39	63.4288	25.7794	0.859728
0.8	2927340	120.161	25.7794	0.859728
0.9	5.4965E+010	434.747	25.7794	0.859728

La figure Fig(7-13) illustre la variation du gain en fonction de la largeur de la bande : le gain de la structure canonique varie sensiblement avec la fréquence de coupure, il est important pour les filtres à bandes étroites. Pour la structure décomposée canonique le gain est relativement moins sensible à la largeur de la bande. Le gain de bruit est insensible à la transformation fréquentielle d'ordre 1 pour les structures optimales et décomposées optimales.

Conclusion

La simulation d'un filtre RII a permis d'étudier les performances des quatre structures d'état : canonique globale, décomposée canonique, décomposée optimale et optimale globales; et de vérifier la correspondance entre la théorie et la pratique .

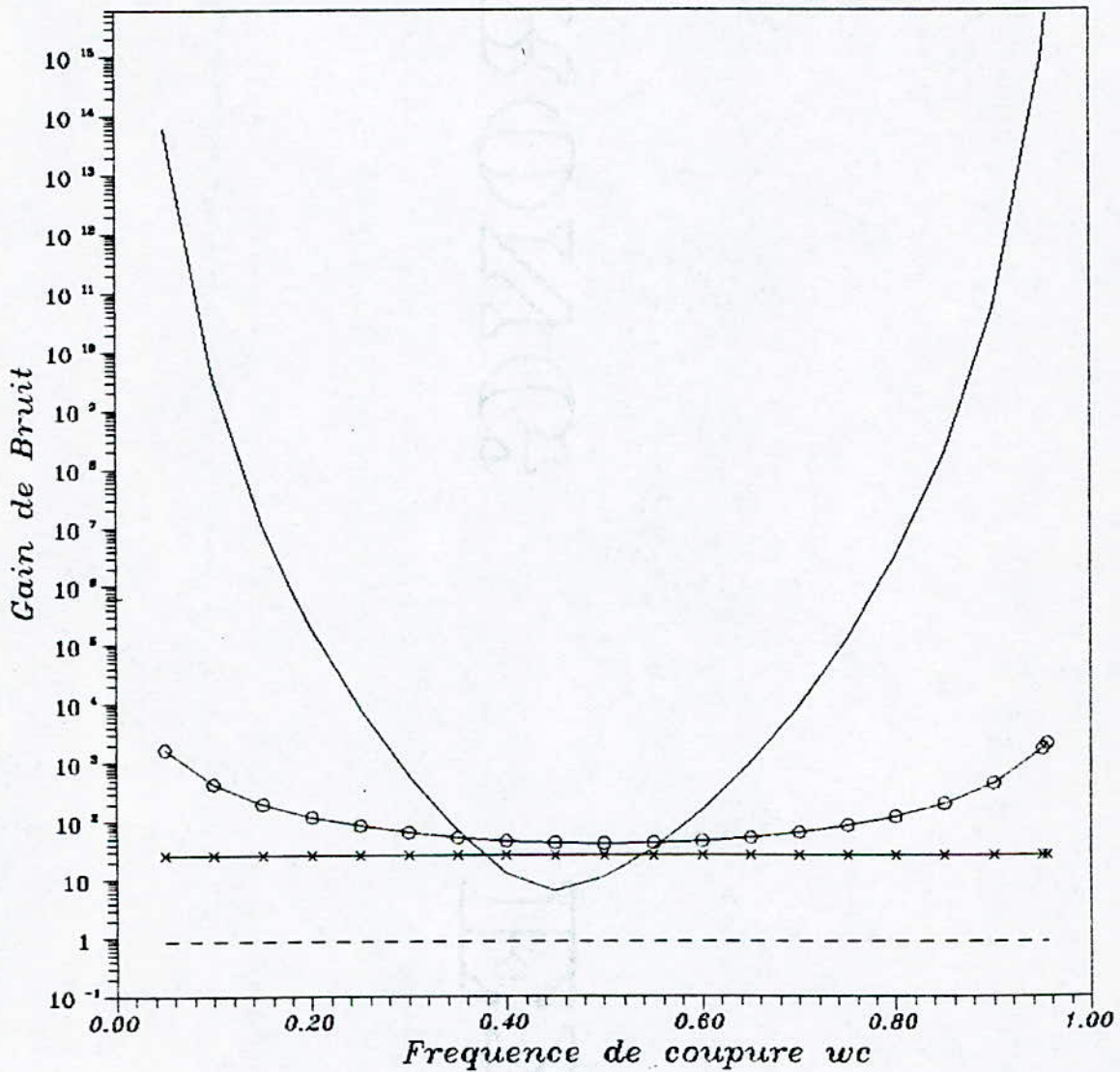


FIG-(7-43) EFFET DE LA TRANSFORMATION FREQUENTIELLE PB-PB SUR LE BRUIT DE CALCUL POUR DIFFERENTES STRUCTURES

- Structure Canonique
- ○ ○ ○ Structure Decomposee Canonique
- x x x x Structure Decomposee Optimale
- - - Structure Optimale

CHAPITRE VIII

CONCLUSION GENERALE

CHAPITRE VIII

CONCLUSION GENERALE

Après avoir exposé brièvement les différentes méthodes de synthèse des filtres numériques ainsi que les distorsions et les effets non linéaires introduits par la quantification des coefficients et des résultats des opérations arithmétiques du filtre numérique RII en virgule fixe, la caractérisation classique de tels filtres par leurs fonctions de transfert ou leurs réponses impulsionnelles s'est révélée insuffisante pour décrire et analyser leurs comportements surtout dans le cas de l'utilisation de mots de longueur finie.

La représentation d'état du filtre est utilisée pour apporter plus de détails et de précision sur le déroulement des opérations et sur la façon avec laquelle s'effectue le filtrage.

Grâce à la propriété de changement de coordonnées dans l'espace d'état par une simple transformation non singulière, des structures éliminant les dépassements, des structures normalisées et des structures à gain de bruit minimale sont possibles pour un même filtre numérique.

Pour trouver les structures à gain de bruit minimal, deux méthodes d'approche ont été utilisées (MULLIS-ROBERTS [16] et HWANG [17]). L'inconvénient de ces structures est leur complexité du point de vue du nombre d'opérations arithmétiques nécessaires dont, en particulier, celui des multiplications qui est de $(N+1)^2$ pour un filtre d'ordre N et ce nombre n'est que $(2N+1)$ pour les structures canoniques.

La théorie de la décomposition du filtre en cellules de second ordre [1],[28] connectées en parallèle a permis d'étudier deux structures de compromis (canonique et optimale) entre le gain de calcul et la complexité des réalisations. Pour la structure décomposée optimale, trois procédures d'optimisation du gain de bruit (BOMAR, BARNES et HWANG [1],[21],[27]) pour les filtres numériques RII du second ordre ont été étudiées et simulées, elles permettent d'obtenir des cellules du second ordre à gain de bruit minimal, ce qui réduit le nombre des multiplications pour le filtre globale à $(2N+1)$. Toutefois, le gain de bruit

cette structure est supérieur à celui de la structure optimale globale.

Le nombre des multiplications pour les structures décomposées canoniques est le même que celui des structures canonique globale, c'est à dire $(2N+1)$, mais leurs gains en bruit sont beaucoup plus faible; ce qui rend ce type de structures intéressant.

La simulation d'un filtre numérique RII du type passe-bas (CHAPITRE VII) a permis de vérifier les résultats élaborés en théorie ainsi que la validité des hypothèses faites sur les erreurs de calcul. Une comparaison entre les deux structures canonique et décomposée optimale a été établie en termes de variance de l'erreur de sortie, de la sensibilité de la fonction de transfert à la quantification de ses coefficients pour chaque structure et de la qualité du filtrage.

En conclusion, la réalisation d'un filtre numérique RII en virgule fixe n'est pas unique; elle dépend de la structure d'état qui doit être choisie de façon à répondre aux spécifications techniques (tolérances en bruit de calcul) et économiques (nombre de multiplieurs nécessaires).

BIBLIOGRAPHIE

- [11] C.T. MULLIS & R.A. ROBERTS, Digital Signal Processing, Addison Wesley, Reading, MA, 1987.
- [12] M.KUNT, Traitement Numérique des Signaux, Dunod 1981.
- [13] J.LIFERMANN, Les méthodes Rapides de Transformation du signal: Fourier, Walsh, Hadamard; Masson 1980.
- [14] J.MAX, Méthodes et Techniques du Traitement du signal, Tome 1, Masson 1987.
- [15] D.ADDOU, " Le filtrage Numérique: Application au Traitement de la Parole, " Projet de Fin D'etudes, ENPA 1987.
- [16] M.LABARRERE et al., Le Filtrage et ses applications, CEPAD 1988.
- [17] M.S.MOAD, "Synthèse des Filtres numériques RII par la méthode des moindres carrés modifiés," Projet de fin d'etudes, ENPA 1990.
- [18] M.BELLANGER, Traitement Numérique du signal, Masson 1987.
- [19] T.M.CLAASEN et al., "Effects of quantization and overflow in recursive digital filters," IEEE trans. Acous, speech & signal processing, Vol.ASSP-26, mars 1979.
- [10] L.R.RABINER & B.GOLD, Theory and Application of Digital Signal Processing, Prentice Hall, Englewood Cliffs, NJ, 1975.
- [11] A.ANTONIOU, Digital Filters:Analysis and Design, Mac Graw Hill, NY, 1979.
- [12] L.B.JACKSON, "Limit Cycles in State-Space Structures for Digital Filters," IEEE Trans. Circuits & Systems, Vol. CAS-26, Jan 1979.
- [13] P.S.DINIZ & A.ANTONIOU, " More Economical State-Space Digital-Filter Structures Which are Free of Constant-Input Limit Cycles," IEEE Trans. Acous, Speech & Signal Processing, Vol. ASSP-34, Aout 1986.
- [14] W.L.MILLS et al., "Digital Filter Realisations without Overflow Oscillations," IEEE Trans. Acous,Speech & Signal Processing,Vol.ASSP-26,Aout 1978.
- [15] C.W.BARNES,"Roundoff Noise and Overflow in Normal Digital Filters," IEEE Trans. Circuits & Systems, Vol.CAS-23, Mars 1979.
- [16] C.T.MULLIS & R.A.ROBERTS,"Synthesis of Minimum Roundoff Noise Fixed Point Digital Filters," IEEE Trans. Circuits & Systems, Sept 1976.

- [17] S.Y.HWANG, " Minimum Uncorrelated Unit Noise in State-Space Digital Filtering," IEEE Trans. Acous, Speech & Signal processing, Vol.ASSP-25, Aout 1977.
- [18] C.T. MULLIS & R.A. ROBERTS," Roundoff Noise in Digital Filters : Frequency Transformations and Invariants," IEEE Trans. Acous., Speech & Signal Processing, Vol.ASSP-24 , Dec 1976.
- [19] B. DERRAS, " New Efficient State-Space Structures for Realisation of Recursive Digital Filters," Thèse de Master, Université du Colorado, 1985.
- [20] B.W.BOMAR, "Minimum Roundoff Noise Digital Filters with Some Power of Two Coefficients," IEEE Trans. Acous., Speech & Signal Processing, Vol.ASSP-34, Oct 1984.
- [21] B.W.BOMAR, " New Second Order State-Space Structures for Realising Low Roundoff Noise Digital Filters," IEEE Trans. Acous., Speech & Signal Processing, Vol.ASSp-33, fev 1985.
- [22] C.W.BARNES, "computationally Efficient Second Order Digital Filter Sections with Low Roundoff Noise Gain," IEEE Trans. Circuits & Systems, Vol. CAS-31, Oct 1984.
- [23] M. ARJMAND & R.A.ROBERTS, "Reduced Multiplier,Low Roundoff Noise Digital Filters," Department of Electrical Engineering,University of Colorado, IEEE 1979.
- [24] C.W.BARNES, "A parametric Approach to the realisation of Second Order Digital Filter Sections," IEEE Trans. Circuits & Systems, Vol. CAS-32, Juin 1985.
- [25] V. TAVSANOGLU & L. THIELE, " Optimal Design of State-Space Digital Filters by Simultaneous Minimisation of Sensitivity and Roundoff Noise," IEEE Trans. Circuits & Systems, Vol. CAS-31, Oct 1984.
- [26] L.B.JACKSON et Al., "Optimal Synthesis of Second Order State-Space Structures for Digital Filters," IEEE Trans. Circuits & Systems, Vol. CAS-26, Mars 1979.
- [27] C.W.BARNES, "On the Design of optimal State-Space Realisations of Second Order Digital Filters," IEEE Trans. Circuits & Systems ,Vol. CAS-31, Juillet 1984.
- [28] C.H.CHEN (ed), Signal Processing Handbook, Marcel DEKKER, INC, NY & BASEL 1988.

STIRRODING

ANNEXE

EXTRA

STI

Le programme STUCMIN effectue le sectionnement du filtre en cellules du second ordre et construit leurs structures canoniques et optimales, leurs gains de bruit et ainsi que leurs matrices K et W respectifs. Les structures optimales sont établies suivant le choix d'une des trois procédures suivantes:

- * Procédure de BOMAR qui est exécutée par la subroutine MINOISE
- * Procédure de BARNES qui est exécutée par la subroutine MINOR.
- * Procédure de HWANG qui est exécutée par la subroutine HWAN.

```

PROGRAM STRUCMIN
REAL*8 P(20),Q(20),Q1(20),D,AA(10,3),BB(10,2),BI(2),CI(2)
Real*8 K(2,2),KI(2,2),WI(2,2),W(2,2),TR(2,2),A(2,2),DEL,G,G1
Real*8 tn(2),td(3)
COMPLEX*16 POL(10,2),p1
Real*8 GM(10),GMN,alp,bet,AI(2,2),B(2),C(2),WK,X,AA1,BB1,CC1
Open(2,file='Natej.dat',Status='NEW')
Open(1,file='Buthba.DAT',status='OLD')
Do 1 i=1,5
  Write(*,*)
1 Continue
  Write(*,*)'      QUEL L''ORDRE DU FILTRE (N=<18) ? '
  Read(*,*) NO
  N=NO+1
  read(1,*) (Q(i),i=1,N)
  read(1,*) (P(i),i=1,N)
  WRITE(2,*)'
  Write(2,*)'      FILTRE D''ORDRE',NO,'
  WRITE(2,*)'
  Write(2,*)
  Write(2,*)'** Les coefficients du numerateurs sont:'
  Write(2,*)
  DO 5 I=1,N
    Write(2,*) Q(i)
5 Continue
  Write(2,*)
  Write(2,*)'** Les coefficients du denumerateur sont:'
  Do 6 I=1,N
    Write(2,*) P(I)
6 Continue
C IN =NOMBRE DE CELLULES DU 2nd ET 1er ORDRE
  IN=NO/2
  if (mod(NO,2).eq.1) Then
    IN=(NO+1)/2
    IMP=1
  Endif

```

C

C DIVISION EUCLIDIENNE DE LA FONCTION DE TRANSFERT

C G= Q/P = d+ Q1/P

D=Q(1)/P(1)

DO 11 I=2,N

Q1(I-1)=Q(I)-P(I)*D

11 CONTINUE

Call decompos(P,Q1,BB,AA,POL,IN,NO)

write(2,*)'***LE COEFF. DU CHEMIN DIRECT ENTREE/SORTIE EST:

write(2,*)' D=',D

Write(*,*)' Quelle le facteur de normalisation (.GE.1.)?'

READ(*,*)DEL

Write(2,*)'

write(2,*)' Decomposition en cellules de structures directes

Write(2,*)'

Write(2,*)

G=0.d0

IR=IN

IF(IMP.eq.1) IR=IN-1

Do 10 i=1,IR

Call Fcan(BB(i,1),BB(I,2),AA(i,2),AA(i,3),A,B,C,K,W)

GM(i)=K(1,1)*W(1,1)+K(2,2)*W(2,2)

Call Sortie(I,A,B,C,K,W,GM(i))

G=G+GM(i)

10 Continue

G1=0.d0

IF(imp.eq.1) Then

x=(1-AA(IN,3)*AA(IN,3))

WK=BB(IN,2)*BB(IN,2)/(x*x)

BB1=BB(IN,2)

AA1=-AA(IN,3)

CC1=1.d0

Call ORDRE1(AA1,BB1,CC1)

G1=WK

Endif

Write(2,*)' *** GAIN TOTAL POUR LA FORME DIRECTE ***'

write(2,12) G

12 Format(5x,' GT =',F16.10)

C Determination de la structure donnant un bruit de calcul minime

C pour chaque cellule du 2nd ordre du filtre.

write(*,*)'

write(*,*)' ***** STRUCTURE OPTIMALE *****

write(*,*)'

write(*,*)'

write(*,*)' (1) Methode de minimisation de BOMAR .

write(*,*)'

write(*,*)' (2) Methode de minimisation de BARNES ,

write(*,*)' (ou structure normale)

write(*,*)'

write(*,*)' (3) Methode de minimisation de HWANG .

write(*,*)'

write(*,*)'

Write(*,*)


```

Write(*,*)
write(*,*)'Quel est votre choix?'
read(*,*) Ichoix
If(Ichoix.eq.1) then
    Write(2,*)
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Write(2,*)
    Do 20 i=1,IR
        Call MINOISE(I,BB(i,1),BB(i,2),AA(i,2),AA(i,3),GM(i),DEL)
20    Continue
    Goto 50
Endif
If(Ichoix.eq.2) then
    write(2,*)
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Write(2,*)
    Write(2,*)
    Do 30 i=1,IR
        Alp=dreal(POL(i,1))
        Bet=dimag(POL(i,1))
        Call MINOR(I,BB(i,1),BB(i,2),alp,bet,GM(i),DEL)
30    continue
    Goto 50
Else
    Write(2,*)
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Write(2,*)'
    Do 40 i=1,IR
        Call Fcan(BB(i,1),BB(I,2),AA(i,2),AA(i,3),AI,BI,CI,KI,WI)
        Call Hwan(KI,WI,K,W,GM(i),TR,DEL)
        Call Trans(TR,AI,BI,CI,A,B,C)
        Call Sortie(i,A,B,C,K,W,GM(I)
40    continue
Endif
50 GMN=0.d0
Do 60 i=1,IR
GMN=GMN+GM(i)
60 Continue
IF(imp.eq.1) THEN
    BB1=BB(IN,2)/DEL
    AA1=-AA(IN,3)
    CC1=DEL
    Call ORDRE1(AA1,BB1,CC1)
Endif

```

STRUCTURE OPTIMALE DE
 *** BOMAR ***

STRUCTURE OPTIMALE DE
 *** BARNES ***

STRUCTURE OPTIMALE DE
 *** HWANG ***

```

Write(2,*)
Write(2,*)'*** GAIN TOTAL DE LA FORME OPTIMALE ***'
Write(2,*)' GT=',GMN
write(*,*)'Les resultats sont dans le fichier'
write(*,*)'***** " NATEJ.DAT" ***** '
Stop
End

```

```

SUBROUTINE MINOISE(I,q1,q2,p1,p2,GMN,DEL)
Real*8 A(2,2),B(2),C(2),q1,q2,p1,p2,W(2,2),K(2,2)
Real*8 v1,v2,v3,v4,v5,v6,v7,v8,GMN,tn(2),td(3)
Real*8 G,DEL,D1
Integer I
D1=1.d0/DEL
v1=q2/q1
v2=dsqrt(v1*v1-p1*v1+p2)
v3=v1-v2
v4=v1+v2
v5=p2-1.d0
v6=p2+1.d0
v7=v5*(v6*v6-p1*p1)
v8=p1*p1/4.d0-p2
A(1,1)=-p1/2.d0
A(2,2)=A(1,1)
B(1)=dsqrt(v7/(2.d0*p1*v3-v6*(1.d0+v3*v3)))*D1
B(2)=dsqrt(v7/(2.d0*p1*v4-v6*(1.d0+v4*v4)))*D1
A(2,1)=dsqrt((B(2)*B(2)+v5)*v8/(B(1)*B(1)+v5))
A(1,2)=v8/A(2,1)
C(1)=q1/(2.d0*B(1))*DEL
C(2)=q1/(2.d0*B(2))*DEL
Call Gauss(A,B,K,0)
CALL Gauss(A,C,W,1)
GMN=K(1,1)*W(1,1)+K(2,2)*W(2,2)
Call Sortie(I,A,B,C,K,W,GMN)
Return
End

```

```

SUBROUTINE MINOR(i,q1,q2,alp,bet,GMN,DEL)
Real*8 A(2,2),A1(2,2),B(2),B1(2),C(2),C1(2),K(2,2)
REAL*8 K1(2,2),W(2,2),W1(2,2),DEL,GF
Real*8 alp,bet,q1,q2,phi,r,v,d1,d2,GMN
Integer I
v=(q2+alp*q1)/bet
phi=datan2(-q1,v)
r=dsqrt(q1*q1+v*v)
A1(1,1)=alp
A1(2,2)=A1(1,1)
A1(2,1)=bet
A1(1,2)=-A1(2,1)
B1(1)=dsin(phi/2.d0)
B1(2)=dcos(phi/2.d0)
C1(1)=r*B1(2)

```



```

C1(2)=r*B1(1)
call Gauss(A1,B1,K1,0)
d1=dsqrt(K1(1,1))
d2=dsqrt(K1(2,2))
A(1,1)=A1(1,1)
A(2,2)=A(1,1)
A(2,1)=bet*d1/d2
A(1,2)=-bet*d2/d1
B(1)=B1(1)/(d1*del)
B(2)=B1(2)/(d2*del)
C(1)=C1(1)*d1*del
C(2)=C1(2)*d2*del
Call Gauss(A,B,K,0)
CALL Gauss(A,C,W,1)
GF=K(1,1)*w(1,1)+k(2,2)*w(2,2)
GNM=2*r*r*d1*d1*d2*d2
Call Sortie(I,A,B,C,K,W,GNM)
Return
End

```

```

Subroutine Huan(K0,W0,K,W,GNM,TR,DEL)
Real*8 K0(2,2),K(2,2),K1(2,2),W0(2,2),W(2,2)
Real*8 TR(2,2),TC(2,2),T(2,2),R(2,2),U1,U2,DEL
Real*8 T1(2,2),V,THETA,GNM,W1(2,2),W2(2,2)

```

C
C
C

TC represente la transformation de CHOLESKY .

```

TC(2,1)=0.d0
TC(2,2)=dsqrt(K0(2,2))
TC(1,2)=K0(1,2)/TC(2,2)
TC(1,1)=dsqrt(K0(1,1)*K0(2,2)-K0(1,2)*K0(1,2))/TC(2,2)
Call Xmat2(TC,W0,W1)
Theta=Datan(1.d0)
If(W1(1,1).ne.W1(2,2)) Then
    Theta=Datan(2.d0*W1(1,2)/((W1(1,1)-W1(2,2)))/2.d0
Endif
R(1,1)=dcos(Theta)
R(2,2)=R(1,1)
R(2,1)=Dsin(Theta)
R(1,2)=-R(2,1)
Call Xmat2(R,W1,W2)
U1=Dsqrt(W2(1,1))
U2=Dsqrt(W2(2,2))
V=U2/U1
K(1,1)=1.d0/(Del*Del)
K(2,2)=K(1,1)
K(1,2)=(U1-U2)/((U2+U1)*K(1,1))
K(2,1)=K(1,2)
W(1,1)=(U1+U2)*(U1+U2)/(4.d0*K(1,1))
W(2,2)=W(1,1)
W(1,2)=(W2(2,2)-W2(1,1))/(4.d0*K(1,1))
W(2,1)=W(1,2)

```

```

      GNM=K(1,1)*W(1,1)+K(2,2)*W(2,2)
C
C   T transformation de normalisation
C
      T(1,1)=dsqrt(1.d0+U)*Del/2.d0
      T(1,2)=-T(1,1)
      U=1.d0/V
      T(2,1)=DSQRT(1.d0+U)*Del/2.d0
      T(1,2)=T(2,1)
C
C   Calcul de la transformation globale TR=TC*R*T
C
      Call Mult2(Tc,R,T1)
      Call Mult2(T1,T,TR)
      return
      End

      Subroutine Trans(TR,AI,BI,CI,A,B,C)
      REAL*8 TR(2,2),Ai(2,2),A(2,2),BI(2,2),B(2,2),CI(2,2),C(2,2)
      REAL*8 TRI(2,2),AH(2,2)
      Call INUMAT2(TR,TRI)
      Call Mult2(TRI,AI,AH)
      Call Mult2(AH,TR,A)
      Call Vecmat(TRI,BI,B,1)
      Call vecmat(TR,CI,C,0)
      Return
      End

      Subroutine INUMAT2(A,B)
      Real*8 A(2,2),B(2,2),D
      D=A(1,1)*A(2,2)-A(1,2)*A(2,1)
      B(1,1)=A(2,2)/D
      B(1,2)=-A(1,2)/D
      B(2,1)=-A(2,1)/D
      B(2,2)=A(1,1)/D
      Return
      End

      Subroutine Xmat2(F,X,T)
      Real*8 X(2,2),F(2,2),T(2,2),S(2,2)
      Do 30 I=1,2
        Do 20 J=1,2
          S(I,J)=0.00
          Do 10 K=1,2
            S(I,J)=S(I,J)+F(K,I)*X(K,J)
10          Continue
20        Continue
30      Continue
      Do 60 I=1,2
        Do 50 J=1,2
          T(I,J)=0.00
          Do 40 K=1,2

```



```

          T(I,J)=T(I,J)+S(I,K)*F(k,j)
40      Continue
50      Continue
60      Continue
      Return
      End

Subroutine Fcan(q1,q2,p1,p2,A,B,C,K,W)
REAL*8 q1,q2,p1,p2,tn(2),td(3),A(2,2),K(2,2),W(2,2)
REAL*8 B(2),C(2)
tn(1)=q1
tn(2)=q2
td(1)=1.d0
td(2)=p1
td(3)=p2
  Call State2(Tn,Td,A,B,C)
  Call Gauss(A,B,K,0)
  Call Gauss(A,C,W,1)
Return
End

Subroutine State2(Q,P,A,B,C)
Real*8 Q(2),P(3),A(2,2),B(2),C(2)
A(1,1)=0.d0
A(1,2)=1.d0
A(2,1)=-P(3)
A(2,2)=-P(2)
B(1)=0.d0
B(2)=1.d0
C(1)=Q(2)
C(2)=Q(1)
Return
End

Subroutine Sortie(i,A,B,C,K,W,GNM)
REAL*8 A(2,2),K(2,2),W(2,2),B(2),C(2),GNM
write(2,*)'          ~~~~~ /
Write(2,*)'          CELLULE', I
write(2,*)'          ~~~~~ /
Write(2,*)'** MATRICE  A'
Write(2,*)
Write(2,*) A(1,1),A(1,2)
Write(2,*) A(2,1),A(2,2)
Write(2,*)
Write(2,*)'** VECTEUR B'
Write(2,*)
Write(2,*) B(1)
Write(2,*) B(2)
Write(2,*)
Write(2,*)'** VECTEUR C'
Write(2,*)
Write(2,*) C(1),C(2)

```

```

Write(2,*)
Write(2,*)'** MATRICE K'
Write(2,*)
Write(2,*) K(1,1),K(1,2)
Write(2,*) K(2,1),K(2,2)
Write(2,*)
Write(2,*)'** MATRICE W'
Write(2,*)
Write(2,*) W(1,1),W(1,2)
Write(2,*) W(2,1),W(2,2)
Write(2,*)
Write(2,*)'** GAIN G',i
Write(2,*)
Write(2,*) GNM
Write(2,*)
c Write(2,*)' G1=',G1
Return
END

```

Subroutine Vecmat(Mat,Vec,Res,Icas)

```

c
c   Produit D'une matrice par un vecteur:
c   Icas= 0   => Vect*Mat
c   Icas= 1   => Mat*Vec
c

```

```

Real*8 Mat(2,2),Vec(2),Res(2),Temp
Do 20 I=1,2
  Res(I)=0.d0
  Do 10 K=1,2
    Temp=Mat(I,K)
    If(Icas.eq.0) Temp=Mat(K,I)
    RES(I)=Res(I)+Vec(K)*Temp
10  Continue
20  Continue
Return
End

```

Subroutine Mult2(A,B,C)

```

Real*8 A(2,2),B(2,2),C(2,2)
Do 10 I=1,2
  Do 20 j=1,2
    C(i,j)=0.d0
    Do 30 k=1,2
      C(i,j)=C(i,j)+A(i,k)*B(k,j)
30  Continue
20  continue
10  Continue
Return
End

```

SUBROUTINE Gauss(F,E,Res,Iflag)

C Sous-programme de résolution d'un système d'équations linéaire


```

C                                     [T][x]=[R]
C   Iflag=0 => Calcul de K
C   Iflag=1 => Calcul de w
C
  Real*8 T(3,3),R(3),Pivot,Temp,Cof
  Real*8 Res(2,2),Temp1,Temp2
  Real*8 F(2,2),E(2)
    Temp1=F(1,2)
    Temp2=F(2,1)
  IF(Iflag.eq.1) Then
    Temp1=F(2,1)
    Temp2=F(1,2)
  Endif
  T(1,1)=F(1,1)*F(1,1)-1.d0
  T(3,3)=F(2,2)*F(2,2)-1.d0
  T(2,2)=Temp1*Temp2+F(1,1)*F(2,2)-1.d0
  T(1,2)=2.d0*F(1,1)*Temp1
  T(1,3)=Temp1*Temp1
  T(2,1)=F(1,1)*Temp2
  T(2,3)=Temp1*F(2,2)
  T(3,1)=Temp2*Temp2
  T(3,2)=2.d0*Temp2*F(2,2)
  R(1)=-E(1)*E(1)
  R(2)=-E(1)*E(2)
  R(3)=-E(2)*E(2)
  N=3
  Ier=0
  Do 50 K=1,N-1
    Pivot=DABS(T(K,K))
    IP=K
    Do 10 I=K+1,N
      If (Pivot.Lt.DABS(T(I,K))) Then
        Pivot=DABS(T(I,K))
        IP=I
      Endif
10    Continue
    If (Pivot.EQ.0.D0) Goto 80
    If (IP.Ne.K) Then
      Do 20 I=K,N
        Temp=T(K,I)
        T(K,I)=T(IP,I)
        T(IP,I)=Temp
20    Continue
      Temp=R(K)
      R(K)=R(IP)
      R(IP)=Temp
    Endif
    Do 40 I=K+1,N
      Cof=T(I,K)/T(K,K)
      Do 30 J=K+1,N
        T(I,J)=T(I,J)-Cof*T(K,J)
30    Continue

```

```

      R(I)=R(I)-Cof*R(K)
40  Continue
50  Continue
   If (T(N,N).Eq.0.D0) Goto 80
   R(N)=R(N)/T(N,N)
   Do 70 I=1,N-1
     K=N-I
     Temp=0.D0
     Do 60 J=K+1,N
       Temp=Temp+T(K,J)*R(J)
60  Continue
     R(K)=(R(K)-Temp)/T(K,K)
70  Continue
     Res(1,1)=R(1)
     Res(1,2)=R(2)
     Res(2,1)=R(2)
     Res(2,2)=R(3)
   Return
80  Ier=1
   write(*,*)'Pivot nul'
   Return
   End

```

```

Subroutine DECOMPOS(A,B1,BB,AA,POL,IN,NO)
REAL*8 A(20),B1(20),D,AA(10,3),BB(10,2)
COMPLEX*16 POL(10,2)

```

```

C
C  CALCUL DES POLES.
C
C    CALL BAIRSTOW(A,NO,IN,POL,AA)
C
C  DECOMPOSITION EN SECTIONS DU 2nd ORDRE.
C
C    CALL RESIDU(POL,B1,NO,IN,BB)
C    CALL CELL(NO,IN,POL,AA,BB)
C    Return
C    End

```

```

SUBROUTINE BAIRSTOW(A,M,IN,POL,AA)

```

```

*-----*
*  SUB. BAIRSTOW RECHERCHE LES COUPLES DE POLES D'UN FILTRE *
*  PAR LA METHODE DE BAIRSTOW. *
*  LES PARAMETRES TRANSMIS SONT: *
*  M   : DEGRES DU POLYNOME. *
*  A(M+1): LES COEFFICIENTS DU DENOMINATEUR DE LA FONCTION *
*  DE TRANSFERT DU FILTRE D'ORDRE M. *
*  IN  : NOMBRE DE CELLULES DU 2nd ORDRE; *
*        IN=M/2   SI M EST PAIR,TOUTES LES CELLULES *
*                SONT DU 2eme ORDRE. *
*        IN=(M-1)/2 SI M EST IMPAIR,LA (IN+1) CELLULE EST *
*                DU 1er ORDRE. *

```



```

* LES RESULTATS DU CALCUL SONT:
* POL(IN,2) = MATRICE DES POLES DONT CHAQUE LIGNE
* REPRESENTE UNE SECTION DU 2nd ORDRE ET
* CONTIENT 2 POLES (qui pourraient etre
* complexes conjuges).
* AA(IN,3) = MATRICE DONT CHAQUE LIGNE REPRESENTE UNE
* SECTION DU 2nd ORDRE ET CONTIENT LES
* COEFFICIENTS DU DENOMINATEUR CORRESPONDANT.
*

```

```

COMPLEX*16 POL(10,2)
REAL*8 A(20),B(20),AA(10,3),R(20),S,S1,P,P1,C,D,X,Y,E
CALL EPSILON(E)
E=E**(2.d0/3.d0)
K=M+1
j=1
100 S=-A(2)/A(1)
P=A(3)/A(1)
C
C INITIALISATION DE LA SOMME S ET DU PRODUIT P DES RACINES.
C
B(1)=A(1)
L=0
60 L=L+1
B(2)=A(2)+S*B(1)
DO 20 I=3,K
B(I)=A(I)+S*B(I-1)-P*B(I-2)
20 CONTINUE
C
C LE VECTEUR B CONTIENT LES COEF. DU POLYNOME QUOTIENT.
C
R(1)=0.D0
R(2)=B(1)
DO 30 I=3,K
R(I)=B(I-1)+S*R(I-1)-P*R(I-2)
30 CONTINUE
C
C LE VECTEUR R CONTIENT LES DERIVEES PARTIELLES DES Bi PAR RAPPORT
C A S ET P.
C
D=R(K)*R(K-2)-R(K-1)*R(K-1)
S1=S-(B(K)*R(K-2)-B(K-1)*R(K-1))/D
P1=P-(B(K)*R(K-1)-B(K-1)*R(K))/D
C
C ITERATION DES VALEURS DE S ET P.
C
IF((ABS(S1-S)+ABS(P1-P)).LT.E) GOTO 50
S=S1
P=P1
GOTO 60
50 D=S1*S1-4*P1
IF(D.LT.0.D0) GOTO 70

```

C
C
C

C SERT SEULEMENT A ORDONNER LES POLES TROUVES.

```
C=-1.DO
IF(S1.GT.0.DO) C=1.DO
D=DSQRT(D)
X=(S1+C*D)/2.DO
Y=P1/X
POL(J,1)=CMPLX(X,0.DO)
POL(J,2)=CMPLX(Y,0.DO)
AA(J,1)=1.DO
AA(J,3)=P1
AA(J,2)=-S1
GOTO 80
70 D=DSQRT(-D)
POL(J,1)=CMPLX(S1/2,D/2)
POL(J,2)=CMPLX(S1/2,-D/2)
AA(J,1)=1.DO
AA(J,3)=P1
AA(J,2)=-S1
80 K= K-2
J=J+1
DO 90 I=1,K
A(I)=B(I)
90 CONTINUE
IF(K.GT.3) GOTO 100
IF(K.EQ.2) GOTO 110
IF(K.EQ.1) GOTO 120
S1=-A(2)/A(1)
P1= A(3)/A(1)
GOTO 50
110 Continue
POL(in,1)=CMPLX(-A(2)/A(1),0.DO)
AA(IN,1)=0.DO
AA(IN,2)=1.DO
AA(IN,3)=A(2)/A(1)
125 FORMAT(2X,4F16.10)
120 Continue
do 121 i=1,in
do 122 j=1,2
write(*,*)'pol(' ,i,' ,',j,')=' ,pol(i,j)
122 continue
121 continue
RETURN
END
```

SUBROUTINE EPSILON(EPS)

```
*-----*
* Cette subroutine calcule le plus petit reel pouvant etre *
* represente par la machine. *
*-----*
```

Real*8 EPS


```

EPS=1.D0
10 If ((1.D0+EPS).GT.1.D0) then
    EPS=EPS/2.D0
    GOTO 10
Endif
EPS=EPS*2.D0
Return
End

SUBROUTINE RESIDU(POL,B1,N0,IN,BB)
Complex*16 pol(10,2),R,Num,Den,v(20)
Real*8 B1(20),BB(10,2)
Integer IN
Do 5 i=1,IN
c
c Calcul de la valeur du polynome numerateur de la fonction de
c transfert par l' algorithme de HORNER . Le resultat est Num.
c
    V(1)=cplx(0.d0,0.d0)
    Do 15 k=1,N0
        V(k+1)=B1(k)+V(k)*pol(i,1)
15 Continue
    Num=v(N0+1)
    Den=cplx(1.d0,0.d0)
    Do 25 j=1,IN
        if(i.eq.j) goto 25
        Den=Den*(pol(i,1)-pol(j,1))
25 Continue
    Do 35 j=1,IN-1
        Den=Den*(pol(i,1)-pol(j,2))
35 Continue
    If((2*IN).EQ.N0) then
        Den=Den*(pol(i,1)-pol(IN,2))
    Endif

c
c R est le residu du 1er pole de la cellule i.
c
    R=Num/Den
    If(DIMAG(POL(IN,1)).EQ.0.d0) goto 65

c
c Les BB sont les coefficients du numerateur (polynome du 1er
c ordre= BB1*z+BB2 )
c
    BB(i,1)= 2*dreal(R)
    BB(i,2)=-2*dreal(pol(i,2)*R)
5 Continue
    Goto 75

c
c Consideration du cas d'un pole reel simple :
c
65 BB(IN,2)=dreal(R)
    BB(IN,1)=0.d0

```

```

75 return
End

SUBROUTINE CELL(NO,IN,POL,AA,BB)
Complex*16 pol(10,2)
Real*8 AA(10,3),BB(10,2)
write(2,*)
Write(2,*)'
Write(2,*)'      ** DECOMPOSITION EN CELLULES DU 2nd ORDRE **
Write(2,*)'
Do 20 i=1,IN
  WRITE(2,*)'
  Write(2,*)' CELLULE',i
  write(2,*)'
  Write(2,*)
  write(2,*) ' ***Les coefficients du numerateur sont : '
  write(2,45) BB(i,1)
  write(2,45) BB(i,2)
  write(2,*) ' ***Les coefficients du denominateur sont : '
  do 55 k=1,3
    write(2,45) AA(i,K)
55 Continue
  IF(I.eq.IN) Then
    If(Mod(NO,2).eq.1) then
      Write(2,*)'      **** UN POLE SIMPLES REEL : '
      Write(2,*) Dreal(POL(IN,1))
    Else
      GOTO 30
    Endif
  Endif
30 write(2,*)' ***Les 2 poles complexes sont : '
  write(2,46) Dreal(pol(i,1)),Dimag(pol(i,1))
45 Format(5x,f16.10)
46 format(5x,f16.10,'+/-j',F16.10)
20 Continue
Return
end

Subroutine ORDRE1(AA1,BB1,CC1)
Real*8 AA1,BB1,CC1
  Write(2,*)' ~~~~~/'
  write(2,*)' CELLULE DU 1er ORDRE '
  Write(2,*)' ~~~~~/'
  Write(2,*)
  WRITE(2,*)'*** A = ',AA1
  Write(2,*)
  WRITE(2,*)'*** B = ',BB1
  Write(2,*)
  Write(2,*)'*** C = ',CC1
  WRITE(2,*)
Return
End

```