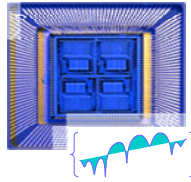


Implantation d'algorithmes de traitement du signal sur les architectures virgule fixe



Daniel Ménard, Olivier Sentieys
 IRISA/ENSSAT Lannion
 Université de Rennes I
sentieys@irisa.fr
<http://www.irisa.fr/R2D2>



UNIVERSITÉ DE RENNES I



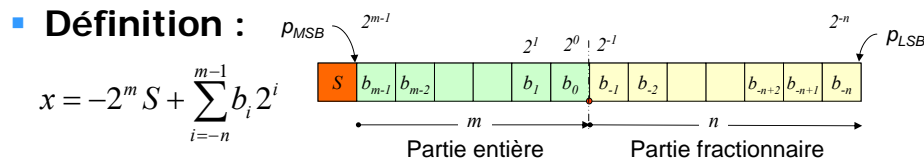
Plan

1. Arithmétique virgule fixe
2. Détermination du domaine de définition
3. Codage des données

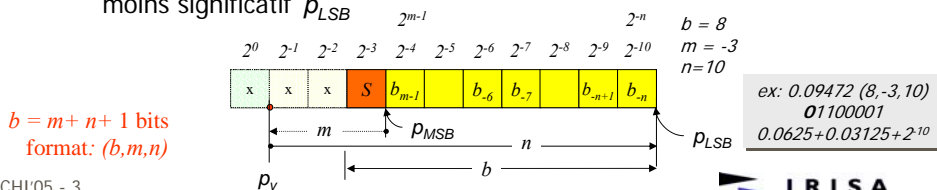
ARCHI'05 - 2



Codage en virgule fixe complément à 2



- m : distance (en nombre de bits) entre la position du bit le plus significatif p_{MSB} et la position de la virgule p_V
- n : distance entre la position de la virgule p_V et la position du bit le moins significatif p_{LSB}



ARCHI'05 - 3

Codage en virgule fixe complément à 2

- Domaine de définition du codage : $[-2^m, 2^m - 2^{-n}]$

- Pas de quantification

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

- Exemple de codage

- $(6, 3, 2)$, Q26

0111.11	7.75
0111.10	7.5
0000.11	0.75
0000.10	0.5
0000.01	0.25
0000.00	0
1111.11	-0.25
1111.10	-0.5
1111.01	-0.75
1000.01	-7.75
1000.00	-8

Annotations:
 -0.5 = -8 + 7.5
 -7.75 = -8 + 0.25



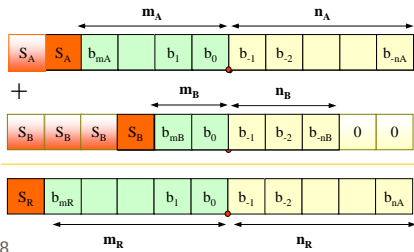
ARCHI'05 - 5

Règles de l'arithmétique virgule fixe

■ Addition: $r = a + b$

- Règle : le format des opérandes a et b doit être identique
- Étapes
 - Choix d'un format commun (b_c, m_c, n_c)
 - Alignement de la virgule
 - Extension des bits des opérandes a et b

$$\begin{cases} m_c = \max(m_A, m_B) \\ n_c = \max(n_A, n_B) \\ b_c = m_c + n_c + 1 \end{cases}$$



$$m_R = \begin{cases} m_c + 1 & \text{si } a + b \notin D_C \\ m_c & \text{si } a + b \in D_C \end{cases}$$

$$n_R = n_c$$

$$b_R = m_R + n_R + 1$$

Exemple de calcul

■ Addition: $r = a + b$

- Débordement possible si $sign(a) = sign(b) = s_{ab}$
 - Débordement présent si $sign(r) \neq s_{ab}$

$$\begin{array}{r} 0110 \text{ (6)} \\ + 1011 \text{ (-5)} \\ \hline 0001 \text{ (1)} \end{array} \quad \begin{array}{r} 0011 \text{ (3)} \\ + 0111 \text{ (7)} \\ \hline 00111 \text{ (7)} \end{array} \Rightarrow \begin{array}{r} 00011 \text{ (3)} \\ + 00111 \text{ (7)} \\ \hline 01010 \text{ (10)} \end{array}$$

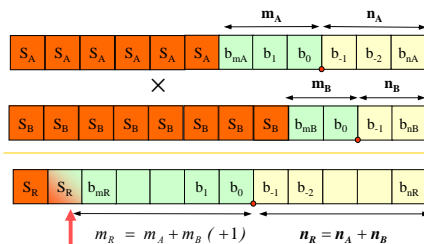
S'il n'y a pas de débordement la dernière retenue peut être ignorée

Nécessité d'un bit supplémentaire pour coder le résultat

Règles de l'arithmétique virgule fixe

■ Multiplication: $r = a \times b$

- Règle : la représentation de a et b doit être identique
- Étapes
 - Extension des bits de signe des opérandes a et b
 - Doublement du bit de signe du résultat
 - Le bit redondant peut être intégré à la partie entière du résultat



$$\begin{cases} n_R = n_A + n_B \\ m_R = m_A + m_B + 1 \\ b_R = b_A + b_B \end{cases}$$

Exemple de calcul

■ Multiplication: $r = a \times b$

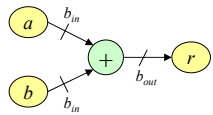
- Extension du signe des opérandes

$$\begin{array}{r} 1100 \text{ (-4)} \\ \times 0110 \text{ (6)} \\ \hline 1111100. \\ 111100. \\ \hline 11101000 \\ -128+64+32+8 = -24 \end{array} \quad \begin{array}{r} 111110 \text{ (-2)} \\ \times 111100 \text{ (-4)} \\ \hline 1110. . \\ 110. . . \\ 10. . . . \\ \hline 001000 \text{ (8)} \end{array}$$

Résumé des règles

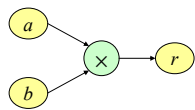
■ Addition

- Choix du format commun



si $b_{out} = b_{in}$ alors $m_c = \max(m_A, m_B, m_R)$
 si $b_{out} > b_{in}$ alors $m_c = \max(m_A, m_B)$

■ Multiplication



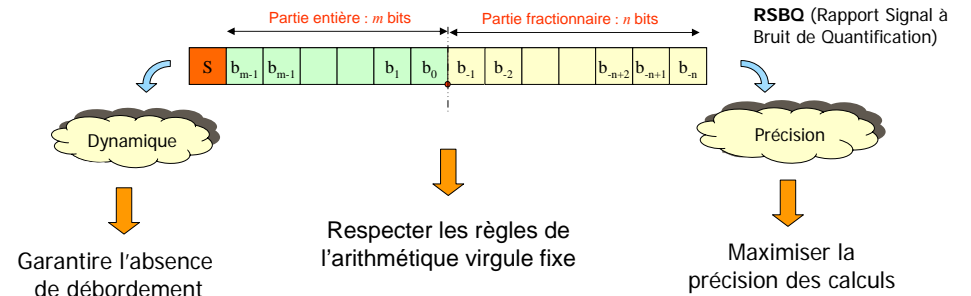
$n_R = n_A + n_B$
 $m_R = m_A + m_B + 1$ (circled) → **Doublement du bit de signe**
 $b_R = b_A + b_B$



Codage en virgule fixe : objectifs

■ Objectifs pour le codage en virgule fixe

- Garantir l'absence de débordements
 - Connaissance du domaine de définition des données
- Maximiser la précision

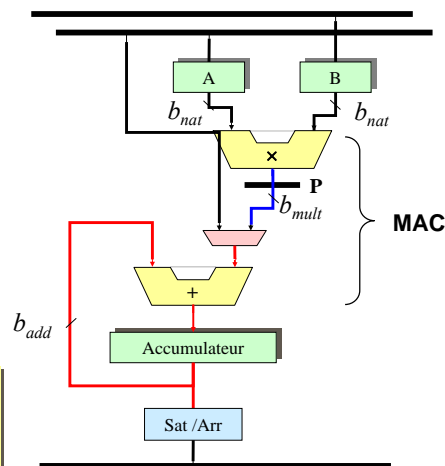


Éléments de l'unité de calcul

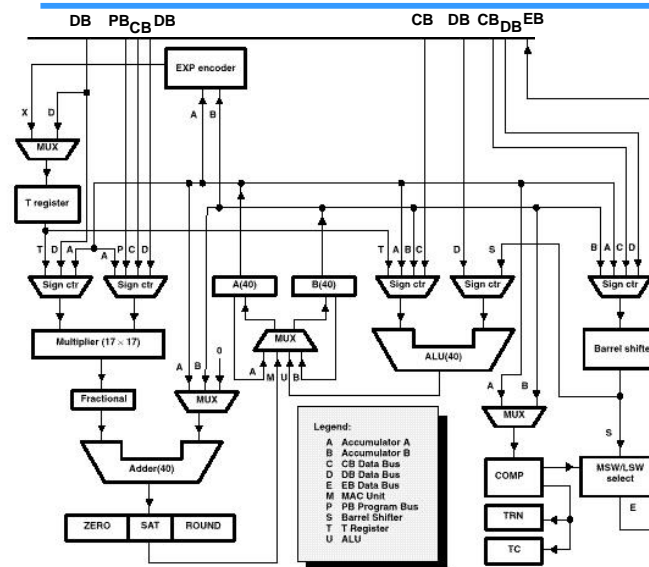
Opérateurs	N _{bits} Entrées	N _{bits} Sortie
Multiplieur	b_{nat}	$2 b_{nat}$
Additionneur/ UAL	$2 b_{nat}$	$2 b_{nat}$
	$2 b_{nat} + b_g$	$2 b_{nat} + b_g$
Saturation/ Arrondi	$2 b_{nat}$	b_{nat}
	$2 b_{nat} + b_g$	b_{nat}

Registre	N _{bits}
Opérande source	b_{nat}
Accumulateur	$2 b_{nat}$
	$2 b_{nat} + b_g$

Interconnexions registres - opérateurs spécialisées
 ⇒ Structure hétérogène
 ⇒ Bonnes performances en terme de consommation et de surface



Exemple : TMS320C54x



1 multiplieur 16*16 bits
 Op source 1 : registre T
 Op source 2 : mémoire
 Op destination :

1 additionneur 40 bits
 1 ALU (40 bits)
 2 registres accumulation 40 bits
 1 registre à décalage en barillet
 1 unité dédiée Viterbi



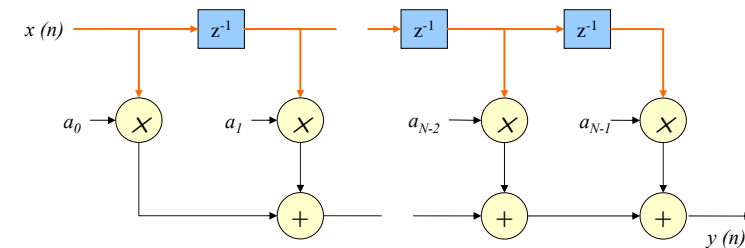
Codage en virgule fixe : étapes

1. Détermination de la dynamique des données : $D(x_i)$
 - Méthodes statistiques ou analytiques
2. Détermination de la position de la virgule : m_{x_i}
3. Détermination de la largeur des données : b_{x_i}
 - Opérateurs SIMD, précision multiple
4. Evaluation de la précision des calculs
 - Méthodes par simulation ou analytiques
 - Calcul du Rapport Signal à Bruit de Quantification



Fil rouge : filtre FIR

$$y(n) = \sum_{i=0}^{N-1} a_i \times x(n-i)$$



Fil rouge : filtre FIR

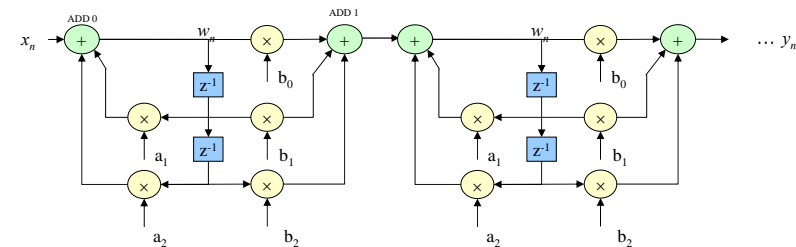
```
float x[N] = {0.123, -0.569, ...} /* Signal d'entrée du filtre */
float h[N] = {-0.0295, 0.0364, ..., -0.0295 }; /* Coefficients du filtre */
int main() {
    float x[N], y[M], acc;
    int i, j;

    for(i=0; i<N; i++){x[i] = 0;} /* Initialisation variables internes */

    for(j=0; j<M; j++) { /* Filtrage du vecteur d'entrée input */
        x[0] = Input[j];
        acc = x[0]*h[0];
        for(i=N-1; i>0; i--)
        {
            acc = acc + x[i]*h[i]; /* Calcul d'une cellule du filtre */
            x[i] = x[i-1]; /* Vieillessement des variables internes */
        }
        y[j] = acc;
    }
}
```

Fil rouge 2 : filtre récursif (IIR)

- Cascade de cellules d'ordre 2



$$w_i(n) = x(n) - \sum_{j=1}^2 a_{j,i} w_i(n-j)$$

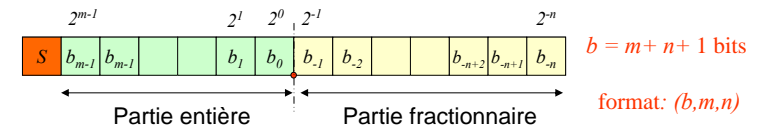
$$y(n) = \sum_{j=0}^2 b_{j,i} w_i(n-j)$$

Plan

- 1. Arithmétique virgule fixe
- 2. Détermination du domaine de définition
- 3. Codage des données

Détermination du domaine de définition

- Format des données: (b, m, n)



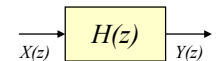
- But: *déterminer le nombre de bits minimal pour représenter la partie entière* (m_{x_i})

$$m_{x_i} = \lceil \log_2(\max(|x_i|)) \rceil$$

Méthodes pour déterminer $D(x_i)$

- Méthodes statistiques
 - Simulation de l'algorithme en virgule flottante
 - Détermination de la dynamique à partir des paramètres statistiques du signal
 - ⇒ Estimation précise mais dépendante du signal d'entrée
- Méthodes analytiques
 - Détermination de l'expression de la dynamique
 - Utilisation de normes (L_1, \dots) ou de l'arithmétique d'intervalle
 - ⇒ Estimation conservatrice mais absence de débordement garantie

Méthodes analytiques



- Normes dans le cas des systèmes linéaires

- Norme L1

$$y_{\max 1} = \max_n(|x(n)|) \sum_{m=-\infty}^{\infty} |h(m)|$$

➤ Méthode garantissant l'absence de débordement

- Norme Chebychev

$$y_{\max 2} = \max_n(|x(n)|) \max_{\omega}(|H(\omega)|)$$

➤ Méthode garantissant l'absence de débordement pour un signal d'entrée à bande étroite ($x(n) = \cos(\omega t)$)

- Norme L2

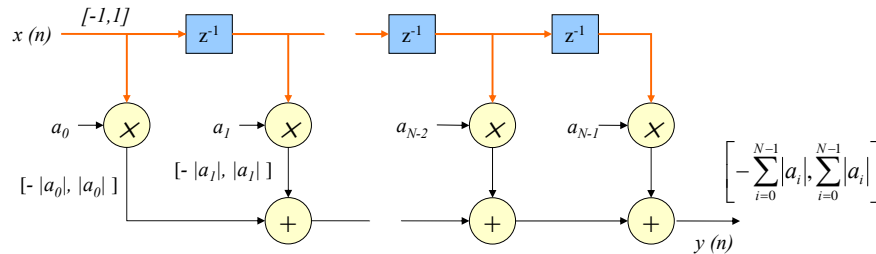
$$y_{\max 3} = \max_n(|x(n)|) \sqrt{\sum_{m=-\infty}^{\infty} |h(m)|^2}$$

➤ Méthode limitant la probabilité de débordement

$$y_{\max 3} \leq y_{\max 2} \leq y_{\max 1}$$

Filtre FIR

- Propagation de la dynamique (intervalles) des entrées au sein du GFS représentant l'application



- Norme L1

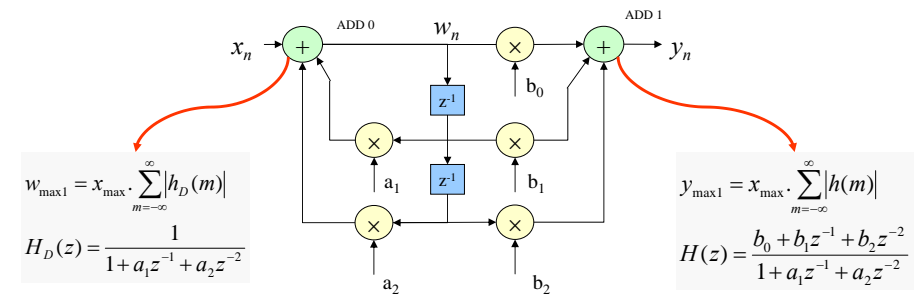
$$\max_n(|y(n)|) = \max_n(|x(n)|) \cdot \sum_{m=-\infty}^{\infty} |h(m)| = \sum_{m=0}^{N-1} |a_m| = 1.65$$

ARCHI'05 - 30



Dynamique des données (IIR)

- Sources de débordements : additionneurs
 - IIR : 2 sources ADD0 et ADD1



$$W_{\max 1} = x_{\max} \cdot \sum_{m=-\infty}^{\infty} |h_D(m)|$$

$$H_D(z) = \frac{1}{1 + a_1 z^{-1} + a_2 z^{-2}}$$

$$y_{\max 1} = x_{\max} \cdot \sum_{m=-\infty}^{\infty} |h(m)|$$

$$H(z) = \frac{b_0 + b_1 z^{-1} + b_2 z^{-2}}{1 + a_1 z^{-1} + a_2 z^{-2}}$$

ARCHI'05 - 31



Dynamique des données (IIR)

- Exemple

$x \in [-1, 1]$

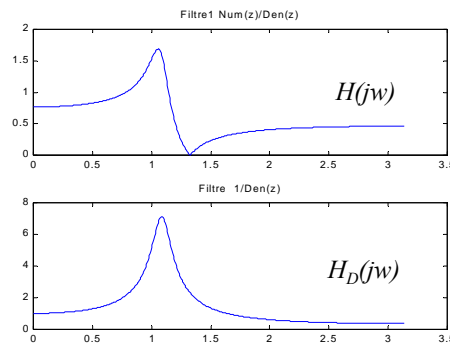
$$H(z) = \frac{0.5 - 0.2428z^{-1} + 0.5}{1 - 0.85z^{-1} + 0.8417z^{-2}}$$

$$x_{\max} \sum_{m=-\infty}^{\infty} |h(m)| = 2.5645 \rightarrow m_y = 2$$

$$x_{\max} \sum_{m=-\infty}^{\infty} |h_D(m)| = 9 \rightarrow m_w = 4 \rightarrow k_1 = 4$$

Codage des coefficients

a1 = -27951
a2 = 27581
b2 = 16384
b0 = 16384
b1 = -7956



Fonction de transfert des filtres $H(z)$ et $H_D(z)$

ARCHI'05 - 32



Dynamique des données (IIR)

- Comparaison des différentes méthodes

Méthodes	$G_{1,max}$	N_1	$G_{2,max}$	N_2
Méthodes analytiques				
Norme L ₁	2,56	2	9	4
Norme Chebychev	1.687	1	7,13	3
Méthodes statistiques				
Chirp	1,66	1	7.12	3
Bruit blanc gaussien	0.74	0	2.34	2
Bruit blanc uniforme	1.4	1	4.87	3

$$G_{1,max} = \frac{y_{\max}}{x_{\max}} \quad N_1 = \lceil \log_2(G_{1,max}) \rceil$$

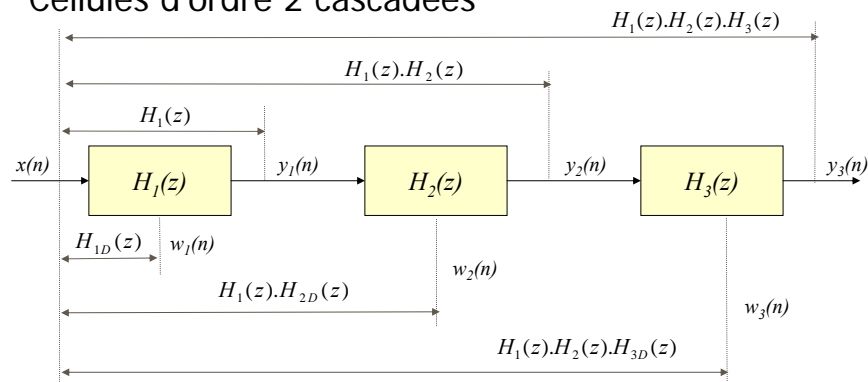
$$G_{2,max} = \frac{W_{\max}}{x_{\max}} \quad N_2 = \lceil \log_2(G_{2,max}) \rceil$$

ARCHI'05 - 33



Dynamique des données (IIR)

- Cellules d'ordre 2 cascadiées



$$H_{iD}(z) = \frac{1}{1 + a_{i1}z^{-1} + a_{i2}z^{-2}} \quad H_i(z) = \frac{b_{0i} + b_{1i}z^{-1} + b_{2i}z^{-2}}{1 + a_{i1}z^{-1} + a_{i2}z^{-2}}$$

Plan

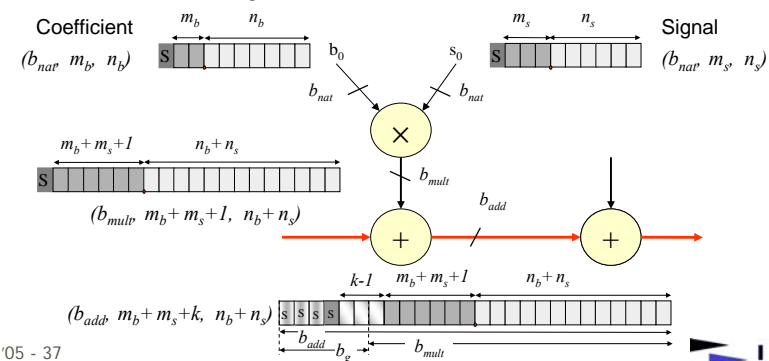
1. Arithmétique virgule fixe
2. Détermination du domaine de définition
3. Position de la virgule et codage des données
Bruits de quantification

Contraintes du codage en virgule fixe

- Respect des règles de l'arithmétique virgule fixe
 - Alignement de la virgule des opérandes sources d'une addition ou d'une soustraction
 - Format commun pour les opérandes sources
- Prévenir les débordements
 - Débordement possible lors d'une addition ou d'une soustraction
 - Nécessité d'introduire des bits supplémentaires
 - Recadrage des données : adapter le format des données à la dynamique des données

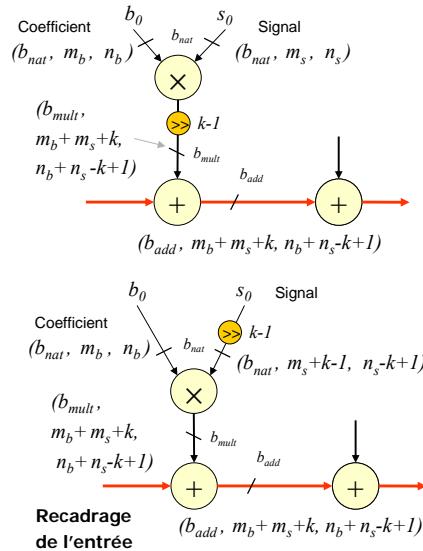
Introduction des bits supplémentaires

- Utilisation des bits de garde de l'accumulateur (b_g bits)
 - $b_{add} = b_{mult} + b_g$
 - Pas de débordement si le nombre de bits supplémentaires k respecte $k \leq b_g$



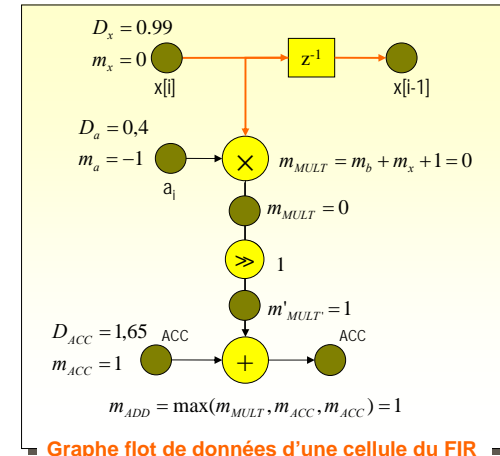
Introduction des bits supplémentaires

- Recadrage de la sortie du multiplieur ($b_{add} = b_{mult}$)
 - Recadrage interne
- Recadrage avant la multiplication ($b_{add} = b_{mult}$)
 - Recadrage de l'entrée
 - Recadrage des coefficients
 - réalisé lors du codage de ces coefficients



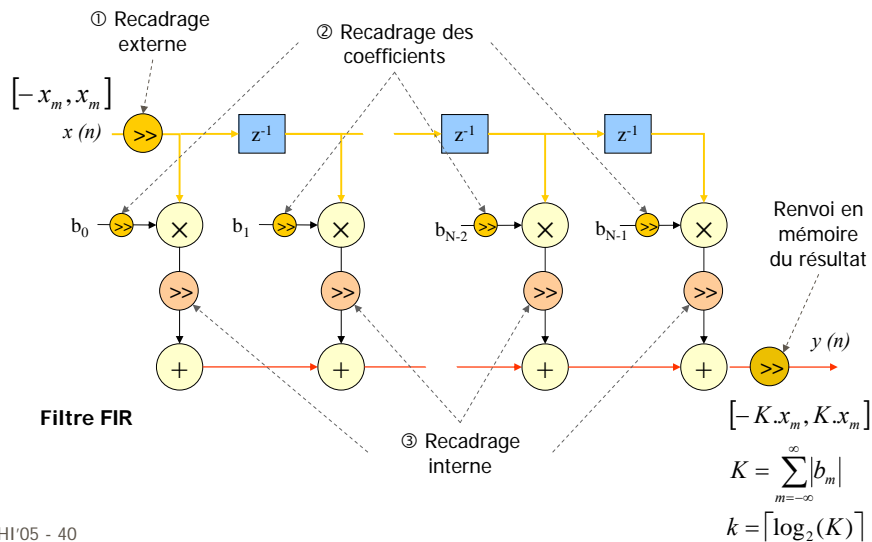
Filtre FIR

- Détermination de la position de la virgule pour les données et les opérations
- Insertion des opérations de recadrage pour aligner la position de la virgule



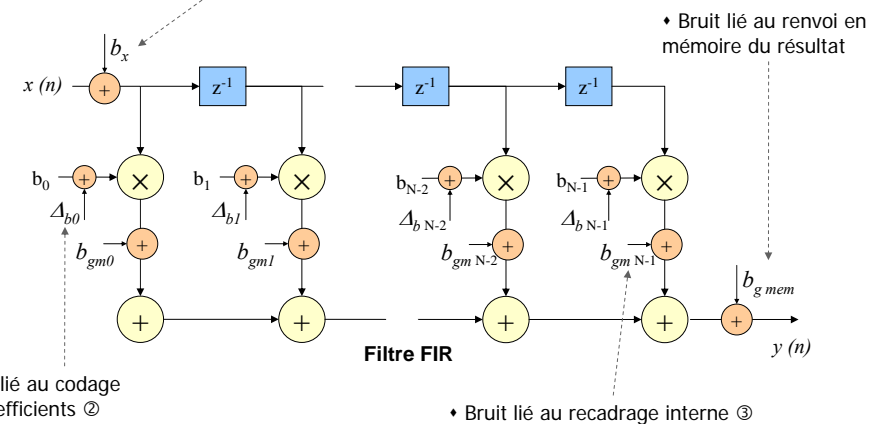
Graphique de données d'une cellule du FIR

Recadrage des données dans un FIR



Sources de bruit dans un FIR

- Bruit de quantification associé à l'entrée
- Bruit lié au recadrage externe ①



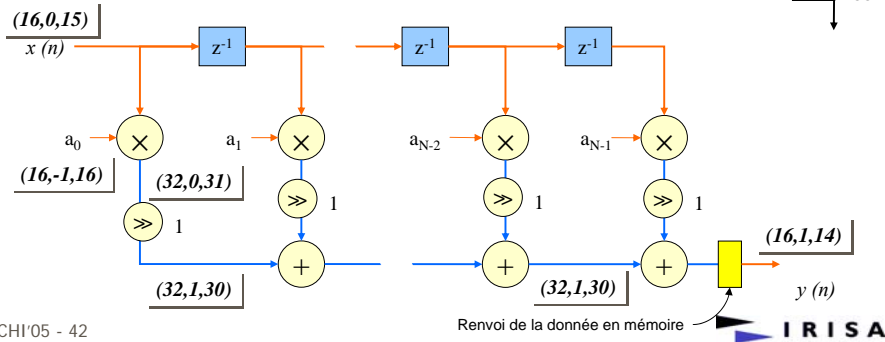
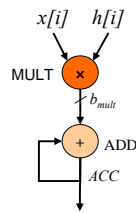
- Biais lié au codage des coefficients ②

- Bruit lié au renvoi en mémoire du résultat

Filtre FIR

Architecture du processeur cible

- Données en mémoire sur 16 bits
- Multiplication 16 bits × 16 bits ⇒ 32 bits
- Addition 32 bits + 32 bits ⇒ 32 bits



Filtre FIR : code C virgule fixe

```

int x[M] = {-809, -6180, -1570, ...} /* Signal x, codage (16,0,15)*/
int h[N] = {-1933, 2386, 3631,}; /* Coefficients (16,-1,16) */

int main() {
    for(i=0; i<N; i++){x[i] = 0;} /* Initialisation des variables internes du filtre */

    for(j=0; j<M; j++){ /* Filtrage du vecteur d'entree input */
        x[0] = Input[j];
        acc = (long)(x[0]*h[0])>>1;

        for(i=N-1; i>0; i--){
            acc = acc + ((long)(x[i]*h[i])>>1); /* Calcul d'une cellule du filtre */
            x[i] = x[i-1]; /* Vieillesse des variables internes */
        }
        y[j] = (int)(acc>>16);
    }
}
    
```

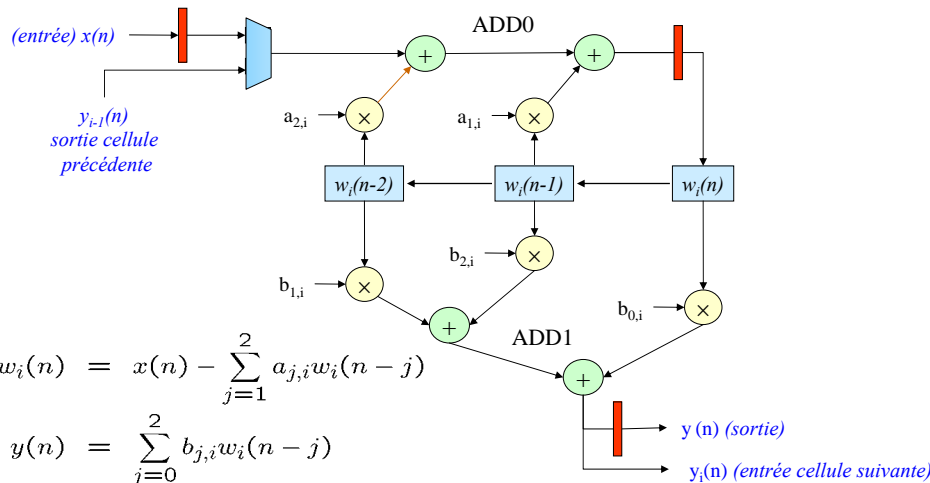
Le signal d'entrée et les coefficients sont spécifiés au niveau du code C en entiers sur 16 bits (absence de type virgule fixe en C) :

- l'entier représentant x (16,0,15) est obtenu en multipliant x par 2^{15}
- l'entier représentant h (16,-1,16) est obtenu en multipliant h par 2^{16}

Recadrage de la sortie de la multiplication : changement de format : (32,0,31) ⇒ (32,1,30)

Réduction de la largeur de la variable 32 bits ⇒ 16 bits
Récupération des 16 bits les plus significatifs de la donnée (l'opération de cast sur `acc` permet de récupérer les bits de poids faible uniquement)

Filtre IIR d'ordre 2 (cellule i)



Dynamique des données (IIR)

Exemple

$x \in [-1,1]$

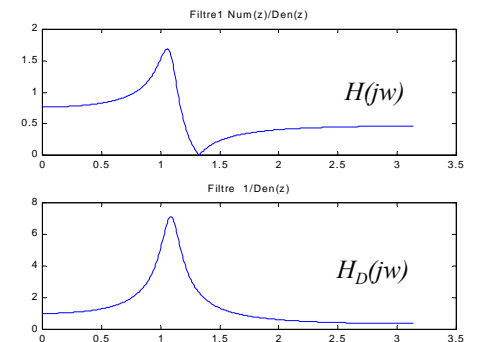
$$H(z) = \frac{0.5 - 0.2428z^{-1} + 0.5}{1 - 0.85z^{-1} + 0.8417z^{-2}}$$

$x_{\max} \sum_{m=-\infty}^{\infty} |h(m)| = 2.5645 \rightarrow m_y = 2$

$x_{\max} \sum_{m=-\infty}^{\infty} |h_D(m)| = 9 \rightarrow m_w = 4 \rightarrow k_1 = 4$

Codage des coefficients

- a1 = -27951
- a2 = 27581
- b2 = 16384
- b0 = 16384
- b1 = -7956



Fonction de transfert des filtres $H(z)$ et $H_D(z)$

Position de la virgule (IIR)

$$m_x = \lceil \log_2(\max_n(|x(n)|)) \rceil = 0$$

$$m_{w_i} = \lceil \log_2(\max_n(|w_i(n)|)) \rceil = 4$$

$$m_{y_i} = \lceil \log_2(\max_n(|y_i(n)|)) \rceil = 2$$

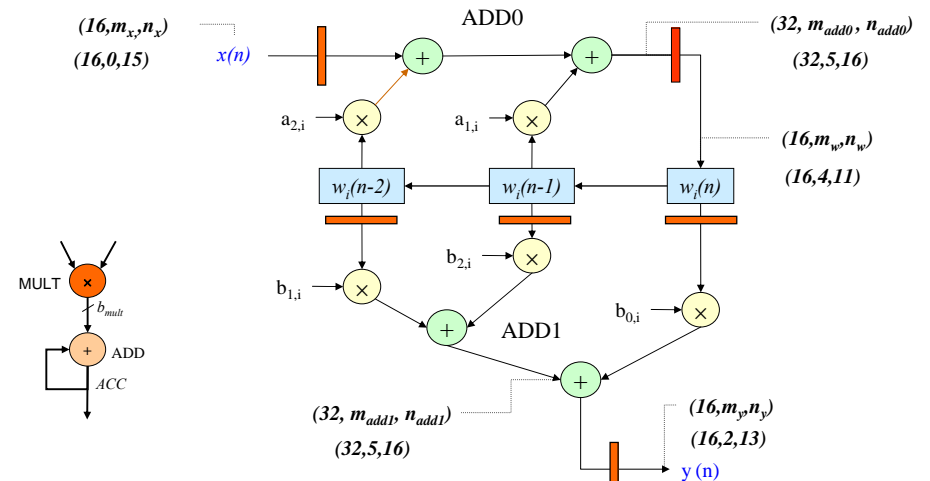
$$m_{a_{ij}} = \lceil \log_2(\max_j(|a_i|)) \rceil = 0$$

$$m_{b_{ij}} = \lceil \log_2(\max_j(|b_{ij}|)) \rceil = 0$$

$$m_{ADD0} = \text{MAX}(m_x, m_{a_i} + m_{w_i} + 1, m_{w_i}) = 5$$

$$m_{ADD1} = \text{MAX}(m_{b_i} + m_{w_i} + 1, m_{y_i}) = 5$$

Codage des données (IIR)



Conclusion

- Différentes étapes du flot de codage
 - Détermination des fonctions de transfert intermédiaires
 - Calcul de leur réponse impulsionnelle
 - Détermination du domaine de définition des données
 - Analyse du graphe flot de données du programme
 - Détermination du système d'équations régissant le format des données
 - Résolution du système d'équations
 - Obtention des recadrages et du format des coefficients

Conversion flottant-fixe (IRISA)

