

Plan

- Circuits logiques
- Représentation des nombres
- Unité Arithmétique et Logique
- Notions de temps et de mémorisation
- Contrôle et jonction des composants
- Evolution des ordinateurs – Historique
- Un microprocesseur simple
- Programmation d'un microprocesseur
- Système complet
- Les microprocesseurs actuels
- Exploitation de la performance des microprocesseurs

Représentation des Nombres

- Deux niveaux de valeur:

V_{DD} et V_{SS} → 1/0.

- Représentation des nombres en base 2.

- Chiffre binaire = *bit* (*binary digit*).

- **Mot n-bits** : nombre de bits utilisés dans les opérateurs d'un processeur.

41: 101001

Bit de poids fort à gauche
(*most significant bit*)

Bit de poids faible à droite
(*least significant bit*)

Nombres Signés

- Trouver une représentation qui affecte peu le coût des circuits de calcul (taille et temps d'exécution) → représentation **transparente** pour principaux opérateurs (addition).

- 1ère solution:

- Signe = bit de poids fort.
- 0: >0, 1: <0.

- Mot = 4 bits.
- Deux représentations de 0:
 - 0000
 - 1000
- Représentation non transparente.

(-1)	1	0	0	1	
(1)	+	0	0	0	1
(-2)	1	0	1	0	

Nombres Signés

• 2ème solution:

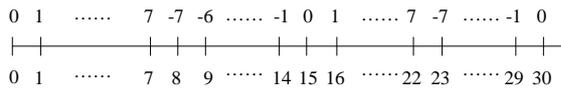
- Opposé = complément
- $-A = A'$.
- $A + (-A) = 0$.
- «Complément à 1».

(-1)	1	1	1	0	(-2)	1	1	0	1		
(1)	+	0	0	0	1	(1)	+	0	0	0	1
(-0)	1	1	1	1	(-1)	1	1	1	0		

- Mot = 4 bits.
- Deux représentations de 0:
 - 0000 (0)
 - 1111 (-0)
- Toujours signe = bit de poids fort.
- Représentation non transparente.

5	1	0	1	0	1
(-2)	+	1	1	0	1
(??)	1	0	0	1	0

Nombres signés



• Complément à 1 = convention:

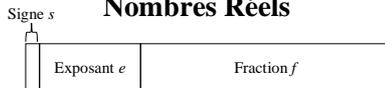
- $0 \leq x \leq 7: X_{CPLT} = x$
- $-7 \leq x \leq 0: X_{CPLT} = x + 15 \Leftrightarrow X_{CPLT} = 15 - |x| = 1111 - x_3x_2x_1x_0 = (1-x_3)(1-x_2)(1-x_1)(1-x_0) = |x|'$
- $x \geq 0, y \geq 0$: si pas de dépassement, $x+y \leq 7$ et $X_{CPLT} + Y_{CPLT} = x+y$.
- $x \geq 0, y \leq 0: X_{CPLT} + Y_{CPLT} = x + y + 15$
 - $x+y \leq 0$: résultat dans $[0;15]$
 - $x+y > 0$: résultat valide mais dans $[16;22]$ (sur 5 bits)
- $x \leq 0, y \leq 0$: si pas de dépassement, $-7 \leq x+y$ et $X_{CPLT} + Y_{CPLT} = 30 + x + y \Rightarrow$ résultat valide mais dans $[23;30]$ (5 bits)

\Rightarrow Résultat addition en complément à 1 obtenu par $Z = (X_{CPLT} + Y_{CPLT}) \text{ modulo } 15$.

Nombres Signés

- Représentation non transparente. Impact sur le circuit ?
 - Revenir sur 4 bits (mot du processeur)
 - Résultat $30 \geq Z > 15 \rightarrow$ sur 5 bits: $1z_3z_2z_1z_0$ et $z_3z_2z_1z_0 \neq 1111$.
 - $Z \text{ mod } 15 = (10000 + z_3z_2z_1z_0) \text{ mod } 15 = (15 + 1 + z_3z_2z_1z_0) \text{ mod } 15 = (1 + z_3z_2z_1z_0) \text{ mod } 15$
 - $\Leftrightarrow 1 + z_3z_2z_1z_0$
- \Rightarrow Addition de la retenue, soit **deux étapes ou deux circuits d'addition**, pour une seule opération d'addition.

Normalisation de la Représentation des Nombres Réels



- Normalisation (IEEE 754)
 - Améliorer la précision/stabilité des résultats et la portabilité des programmes.
 - Représentation à la fois des très grands nombres et des très petits nombres.

Nombre de bits	Taille de f	Taille de e	E_{min}	E_{max}
32 Simple précision	23	8	-126	+127
64 Double précision	52	11	-1022	+1023

Norme IEEE 754

- Exposant:
 - **Biaisé:** e
 - **Non biaisé:** $e - E_{max}$
 - Intérêt: ordre lexicographique \Rightarrow comparaison facilitée. Sur 4 bits:
 - biais = 7; $-1 \rightarrow 0110$, $1 \rightarrow 1000$
 - complément à 2: $-1 \rightarrow 1111$, $1 \rightarrow 0001$
- Ordre signe-mantisse-exposant:
 - Ordre lexicographique pour nombres positifs \Rightarrow comparaison facilitée.
 - Exemple (f sur 4 bits, e sur 4 bits):
 - $1.10x2^2 \rightarrow 0\ 1001\ 1000$
 - $1.01x2^1 \rightarrow 0\ 1011\ 0100$
- $E_{min} \leq e - E_{max} \leq E_{max}$:
 - 1 «implicite» pour que mantisse ≥ 1 .
 - nombres **normalisés**.

$e - E_{max}$	f	Nombre représenté
$E_{min} - 1$	0	± 0
$E_{min} - 1$	$\neq 0$	$(-1)^s x 0.f x 2^{E_{min}}$
$E_{min} \leq e - E_{max} \leq E_{max}$	-	$(-1)^s x 1.f x 2^{e - E_{max}}$
$E_{max} + 1$	0	$\pm \infty$
$E_{max} + 1$	$\neq 0$	NaN

- 2 représentations de 0:
 - 0^+ ($s=0$) et 0^- ($s=1$).
 - 0^+ : tous les bits à 0
 - La norme impose que le test $0^+ = 0^-$ soit vrai.
- Infini:
 - $+\infty$ ($s=0$)
 - $-\infty$ ($s=1$)

Norme IEEE 754

- $e = E_{min} - 1, f \neq 0$:
 - Nombres **dénormalisés**.
 - Nombres inférieurs à $1.0 x 2^{E_{min}}$ peuvent être représentés
 - **Underflow** graduel: limiter les erreurs d'arrondi avec les petits nombres.
 - Exemple (Mantisse = 4 bits):
 - $A = 1.0010 x 2^{E_{min}}$, $B = 1.0001 x 2^{E_{min}}$
 - Normalement $A - B$ résultat non représenté, donc 0, or $A \neq B$.
 - Avec nombres dénormalisés: $A - B = 0.0001 x 2^{E_{min}}$
 - Exemple:
 - Explorer l'espace des valeurs d'une fonction f dont les caractéristiques sont inconnues.

$e - E_{max}$	f	Nombre représenté
$E_{min} - 1$	0	± 0
$E_{min} - 1$	$\neq 0$	$(-1)^s x 0.f x 2^{E_{min}}$
$E_{min} \leq e - E_{max} \leq E_{max}$	-	$(-1)^s x 1.f x 2^{e - E_{max}}$
$E_{max} + 1$	0	$\pm \infty$
$E_{max} + 1$	$\neq 0$	NaN

- **NaN (Not a Number):**
 - $\sqrt{-1}, \frac{0}{0}, \infty - \infty$
 - Evite une interruption du programme (doit être gérée et dépend de l'environnement).
 - Si une des variables est **NaN**, la sortie est toujours **NaN**.

Exceptions

- Par défaut, continuer dans arithmétique étendue:
 - $\frac{1}{0^+} = +\infty$
 - $\frac{1}{-\infty} = 0^-$
 - $(+\infty) + (+\infty) = +\infty$
 - $+\infty \times (-\infty) = -\infty$
 - $+\infty - \infty = NaN$
 - $\frac{0}{0} = NaN$
 - $0 \times \infty = NaN$
 - $F(\dots, NaN, \dots) = NaN$
 - Dépassement de capacité $\rightarrow NaN$
- En cas d'exception, choix possible entre:
 - Interrompt le programme (gestion exception par l'environnement logiciel).
 - L'utilisateur gère l'erreur dans le programme (des drapeaux signalent les exceptions).
 - Continuer le calcul.
- Exceptions:

Drapeau	Résultat
<i>Underflow</i>	$0, \pm x_{\min}$, nombre dénormalisé
<i>Overflow</i>	$\pm \infty, x_{\max}$
Division par 0	$\pm \infty$
Résultat NaN et opérandes $\neq NaN$	NaN
Inexact	Résultat arrondi

Quelques Exemples

Normalisé:

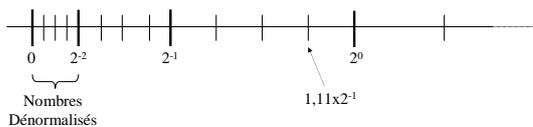
$$s, f, e \rightarrow (-1)^s \times \left(1 + \sum_{i=1}^{23} f_i 2^{-i}\right) \times 2^{e-E_{\max}}$$

$e-E_{\max}$	f	Nombre représenté
$E_{\min}-1$	0	± 0
$E_{\min}-1$	$\neq 0$	$(-1)^s \times 0, f \times 2^{E_{\min}}$
$E_{\min} \leq e \leq E_{\max}$	-	$(-1)^s \times 1, f \times 2^{e-E_{\max}}$
$E_{\max}+1$	0	$\pm \infty$
$E_{\max}+1$	$\neq 0$	NaN

- 32 bits.
- A= 1 10000010 001100000000000000000000
 - $s=1$ (-), $e=130-127=3$, $f=2^{-3}+2^{-4}=0,125+0,0625$
 - A= $-1,1875 \times 2^3 = -9,5$
- A= 1 11111111 001100000000000000000000
 - $e=255 \rightarrow e-E_{\max}=E_{\max}+1, f \neq 0$
 - A= NaN
- A= 1 00000000 001100000000000000000000
 - $s=1$ (-), $e=0 \rightarrow e-E_{\max}=E_{\min}-1, f=0,1875$
 - A= $-0,1875 \times 2^{-126}$ (dénormalisé).

Arrondis

- Pour chaque valeur d'exposant, distance double entre valeurs de mantisse.
- Mais, erreur relative constante.
- Exemple: $f=2$ bits, $e=3$ bits.



Arrondis

- 4 types d'arrondi fournis par la norme:
 - Vers 0
 - Vers $+\infty$
 - Vers $-\infty$
 - Au plus près
- Théoriquement, on peut sélectionner le mode d'arrondi, mais choix souvent imposé dans les langages.
- Arrondi **exact** sur les opérations élémentaires (+, -, ×, /, √):
 - Résultat = arrondi (calcul effectué avec une précision infinie).
 - Pour obtenir ce résultat, modes étendus non accessibles à l'utilisateur:
 - 32 bits, simple précision: $e \geq 11, f \geq 32, \geq 6$ octets.
 - 64 bits, double précision: $e \geq 64, f \geq 15, \geq 10$ octets.
- L'arrondi exact n'est pas encore garanti pour les fonctions élémentaires (*sin, cos, ...*)

Multiplication Flottante

- $f=3$ bits, $e=4$ bits.
- -7 x 15:
 - $x_1 = -7 \rightarrow -1,110 \times 2^2 \rightarrow 1\ 1001\ 110$
 - $x_2 = 15 \rightarrow 1,111 \times 2^3 \rightarrow 0\ 1010\ 111$
- $s = s_1 \oplus s_2$.
- Addition des exposants:
 - $e_1 + e_2 - E_{max}$: chaque exposant contient un biais; après addition, il faut donc en soustraire un biais ($-E_{max}$).
 - $\rightarrow e = 1100$ soit 5 (biais=7 sur 4 bits; addition sur 5 bits).
- Multiplication entière des mantisses (en parallèle):
 - Décompacter f_1 et f_2 (bit 1 implicite):
 - $f_1^p = 1,110, f_2^p = 1,111$
 - $f^p = 11,010010$ sur 8 bits.
 - Normaliser:
 - Revenir à 1,....
 - Incrémenter exposant de 1.
 - $f^p = 1,1010010$
 - Arrondir (au plus près):
 - Résultat: $1, b_1 b_2 b_3$
 - $11010010 \rightarrow$ bits d'arrondi
 - $f = 101$
- Résultat:
 - $1\ 1101\ 101$
 - Soit -104 (perte de précision)

Précision des Calculs

- La norme n'apporte pas de garantie sur la précision de séquences de calcul.

$$F = \frac{x^2}{\sqrt{1+x^3}}$$

- $x \rightarrow +\infty, F \rightarrow +\infty$.
- Pour x grand, x^3 devient $+\infty$, donc F devient 0.
- Lorsque x^2 devient $+\infty$, F devient NaN .
- Pas de distinction entre ∞ et dépassement de capacité.

- Un test dont une opérande est NaN est faux sauf pour l'égalité.
- $\Rightarrow a \geq b$ et $a < b$ ne sont pas exactement contraires.

$$(-1,0111 \times 2^4 + 1,0111 \times 2^4) + 1,1000 \times 2^{-3}$$

$$= 1,1000 \times 2^{-3}$$

$$-1,0111 \times 2^4 + (1,0111 \times 2^4 + 1,1000 \times 2^{-3})$$

$$= -1,0111 \times 2^4 + 1,0111 \times 2^4 = 0$$

- L'addition n'est pas associative.
- Absorption après décalage.

Interprétation des Erreurs d'Arrondi

- $i=100000$
 - $S_{i,i} = 0\ 10000010\ 10000010111010000010110$
 - $1/i = 0\ 01101110\ 01001111100010110101100$
- Aligner les exposants:
 - $10000010-01101110 = 20$; décalage de 20 bits à droite de la mantisse
 - Addition + arrondi

1,10000010111010000010110
+0,000000000000000000001010
1,1000001011101000100000
- $i=2100000$
 - $S_{i,i} = 0\ 10000010\ 11101100111010101111100$ (15,403683)
 - $1/i = 0\ 01101001\ 1111111010011100011111$
- Aligner les exposants:
 - $10000010-01101001 = 25$; décalage de 25 bits à droite de la mantisse
 - Addition + arrondi

1,1101100111010101111100
+0,000000000000000000000000
1,1101100111010101111100
