

## Arithmétique entière des ordinateurs (représentation)

Écriture décimale : écriture positionnelle.

$$\text{Ex : } 128 = 1 \times 10^2 + 2 \times 10^1 + 8 \times 10^0$$

Circuit en logique binaire  $\Rightarrow$  Écriture binaire (base 2)

$$\text{Ex : } (101)_2 = 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0$$

En particulier : Entiers = mots sur l'alphabet  $\{0, 1\}$   
(les chiffres binaires sont appelés bits pour **binary digits**)

En pratique : Entiers mots de longueur fixe ( $T$ )

Selon les machines/constructeurs :  $T$  vaut 16, 32 voire 64 ou même 128.

Un entier « machine » est un mot :

$$b_{T-1}b_{T-2} \dots b_1b_0$$

où  $b_i \in \{0, 1\}$

$$\text{Ex : } (0000001000000001)_2 = 1 \times 2^9 + 1 \times 2^0$$

Note : il y a  $2^N$  mots différents de longueur  $N$  sur un alphabet à deux éléments  $\Rightarrow 2^N$  entiers différents représentables (lesquels ?)

## Arithmétique entière des ordinateurs (opérations)

Table d'addition :

+	$(0)_2$	$(1)_2$
$(0)_2$	$(00)_2$	$(01)_2$
$(1)_2$	$(01)_2$	$(10)_2$

Ex :

$$\begin{array}{r}
 (1000010000000001)_2 \\
 + (1000000000000001)_2 \\
 \hline
 (10000010000000010)_2
 \end{array}$$

Table de multiplication :

×	$(0)_2$	$(1)_2$
$(0)_2$	$(0)_2$	$(0)_2$
$(1)_2$	$(0)_2$	$(1)_2$

Ex :

$$\begin{array}{r}
 (0011)_2 \\
 \times (0010)_2 \\
 \hline
 (0000)_2 \\
 + (0011.)_2 \\
 + (0000..) _2 \\
 + (0000...) _2 \\
 \hline
 (0000110)_2
 \end{array}$$

Problème : débordement de la représentation ou dépassement de capacité

## Arithmétique entière des ordinateurs (autres bases)

Binaire : base 2

Deux chiffres : 0 et 1

Octal : base 8

Huit chiffres : 0, 1, 2, 3, 4, 5, 6 et 7

Héxadécimal : base 16

Seize chiffres : 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E et F

Conversions par « re(dé)groupement »

Ex :  $(101)_2 = (5)_8 = (0101)_2 = (5)_{16}$

$(111101010000)_2 = (7520)_8 = (F50)_{16}$

### Exercices

- Soustraire deux nombres en base 2
- Additionner deux nombres en base 16
- Multiplier deux nombres en base 8
- Convertir un nombre en base 16 en base 8
- Écrire en pseudo-langage/C-- une fonction qui prend un tableau de caractères représentant un nombre écrit en base 8 et renvoie la valeur de type Entier correspondante :  $\{'7', '3'\} \rightarrow 59$

## Exercices (suite)

- Écrire en pseudo-langage/C-- des fonctions de conversions entre les différentes bases (2, 8 et 16). Les nombres étant représentés par des tableaux de caractères.
- Quels sont les mots sur 0,1 représentant des nombres multiples de 2, de 4, de  $2^p$ . Écrire une fonction (pseudo/C--) qui renvoie vraie si le nombre représenté comme tableau de caractères est un multiple de  $2^p$ . La fin du nombre est marquée par le caractère d'espace.
- En base 16, comment sont les mots représentant des multiples de 16 ?

## Arithmétique entière des ordinateurs (nombres négatifs)

Écriture décimale signée : écriture positionnelle + signe

$$\text{Ex : } -128 = -1 \times (1 \times 10^2 + 2 \times 10^1 + 8 \times 10^0)$$

Problème :  $-0$  et  $+0$ , deux mots différents pour le même nombre !

Pas de symbole supplémentaire en machine !  $\Rightarrow$  Un choix parmi de nombreuses possibilités, dont :

- utilisation d'une des lettres du mot pour le signe, il n'y a que 2 signes. Problème : deux représentations  $\neq$  pour 0 ;
- utilisation d'une base exotique (négative ou complexe). Problème : c'est exotique ;
- représentation complémentée. Problème :  $\neq$  selon la complémentation choisie.

En pratique : complément à 2.

## Arithmétique entière des ordinateurs (complément à 2)

Pas de signe explicite !

Nombre de longueur  $n \Rightarrow$  calculs modulo  $2^n$

$$x + (-x) = 0 \text{ mod } 2^n$$

Ex :

$$(00000000)_2 + (00000000)_2 = (00000000)_2 = 0 \text{ mod } 2^8$$

$$(00000001)_2 + (11111111)_2 = (100000000)_2 = 0 \text{ mod } 2^8$$

De façon arbitraire :

chiffre à gauche = 0, alors nombre positif, négatif sinon

$$\text{Ex : } (11111111)_2 = -1, (01111111)_2 = 127$$

L'opération d'addition s'effectue comme sur les nombres non signés

### Exercices

Sur les nombres de 8 bits, et en complément à 2 :

- Calculer les inverses de 35, -89, 127 et -128
- Calculer les sommes  $(-3)+(-7)$ ,  $(-127)+(-1)$ ,  $(127)+(1)$ .

Pourquoi certaines machines permettent-elles de savoir si il y a eu un débordement sur l'avant dernier bit ?

- Comment calculer une soustraction ?

## Arithmétique réelle des ordinateurs (représentation)

Basée sur la notation scientifique

Ex : Nombre d'Avogadro  $\simeq 6.022 \times 10^{23}$

Il suffit d'enregistrer la mantisse et l'exposant

Sur machine on utilise très souvent la norme IEEE 754

Un nombre réel flottant est un triplet  $(s, e, m)$  où :

- $s$  désigne le signe (0 si positif, 1 si négatif) ;
- $e$  est l'exposant par excès de la base ;
- $m$  est la mantisse normalisée.

L'excès ( $q$ ) et la base ( $b$ ) sont implicitement connus ainsi :

$$(s, e, m) = -1^s \times m \times b^{(e-q)}$$

Le mot  $(s, e, m)$  est de taille fixe

Pour IEEE 754 on a :

	$ s $	$ e $	$ m $	q	Taille totale
Simple Précision	1	8	23	127	32
Double Précision	1	11	52	1023	64

## Arithmétique réelle des ordinateurs (représentation)

Pour la simple précision, le mot

$$se_7 \dots e_0 m_{22} m_{21} \dots m_1 m_0$$

représente le nombre

$$-1^s \times 1, m_{22} m_{21} \dots m_1 m_0 \times 2^{e-127}$$

On dit que la mantisse est normalisée

Ex : Le nombre 0,25 s'écrit  $(0,01)_2$  soit

00111110100000000000000000000000

Note : 0 n'est pas représentable !!! Diverses conventions permettant de représenter 0 et les « infinis ».

## Arithmétique réelle des ordinateurs (opérations)

L'addition consiste en :

1. l'alignement des exposants (le plus petit sur le plus grand) ;
2. addition des mantisses ;
3. puis renormalisation.

Attention : l'alignement peut conduire à perdre des bits de la mantisse.

## Arithmétique réelle des ordinateurs (problèmes)

On ne représente pas tous les réels ( $\mathbb{R}$  n'est pas énumérable)

La distribution des réels « machine » sur  $\mathbb{R}$  n'est pas uniforme

Les opérations ne sont pas exactes :  $(a + b) + c$  n'est pas toujours égal à  $a + (b + c)$

Tester l'égalité de deux réels est dangereux

La distance entre deux réels représentables n'est pas toujours représentable

### Exercices

- Représenter en simple précision les nombres  $\frac{1}{32}$ ,  $\frac{1}{3}$  (on sait que  $\frac{1}{3} = \sum_{n=1}^{\infty} \frac{1}{4^n}$ ) et  $2^{30}$ .
- Calculer  $\frac{1}{2} + \frac{1}{8}$ ,  $2^{30} + 1$ .
- Que penser du fragment de programme suivant :

```
Variable a : réel;
```

```
Pour a allant de 2^30 à 2^31 par pas de 1 faire quelque chose
```

- Combien de réels peuvent être représentés en simple (double) précision ?
- Quel est le plus petit (grand) réel représentable ?