



Architecture des ordinateurs

Département Informatique

Erwan LEBAILLY — Vilavane LY — Vincent TRÉLAT — Benjamin ZHU

28 février 2022

Table des matières

1	Chapitre 1	2
1.1	Exercice 1	2
1.2	Exercice 2	2
1.3	Exercice 3	2
1.4	Exercice 4	3
1.5	Exercice 5	4
1.6	Exercice 6	4
1.7	Exercice 7	4
1.8	Exercice 8	4
1.9	Exercice 9	4
1.10	Exercice 10	5
1.11	Exercice 11	5
1.12	Exercice 12	7
1.13	Exercice 13	7
2	Chapitre 2	8
2.1	Exercice 1	8
2.2	Exercice 2	9
2.3	Exercice 3	9
2.4	Exercice 4	9
2.5	Exercice 5	9
2.6	Exercice 6	10
2.7	Exercice 7	10
2.8	Exercice 8	10

1 Chapitre 1

1.1 Exercice 1

Avec la convention $0 \leftrightarrow \text{faux}$ et $1 \leftrightarrow \text{vrai}$, $0 \wedge 1 = \text{faux}$.

1.2 Exercice 2

On donne la table de c_0 :

$$c_0:$$

$a_0 \backslash b_0$	0	1
0	0	1
1	1	0

On peut interpréter cette table comme la table de vérité du "ou exclusif", le *xor*. Ainsi, c_0 coïncide avec $a_0 \oplus b_0 = (a_0 \vee b_0) \wedge (\neg(a_0 \wedge b_0))$.

1.3 Exercice 3

a. Montrons que l'opérateur \oplus est associatif et commutatif :
Soient $a, b, c \in \{0, 1\}$.

Associativité : on donne ci-dessous la table de vérité de $(a \oplus b) \oplus c$ et $a \oplus (b \oplus c)$:

a	b	c	$a \oplus b$	$(a \oplus b) \oplus c$	$a \oplus (b \oplus c)$
0	0	0	0	0	0
0	0	1	0	1	1
1	0	0	1	0	0
1	0	1	1	0	0
0	1	0	1	1	1
0	1	1	1	0	0
1	1	0	0	0	0
1	1	1	0	1	1

Commutativité : on donne ci-dessous la table de vérité de $a \oplus b$ et $b \oplus a$:

a	b	$a \oplus b$	$b \oplus a$
0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	0

Enfin, $a \oplus a = 0$ et $a \oplus 0 = a$.

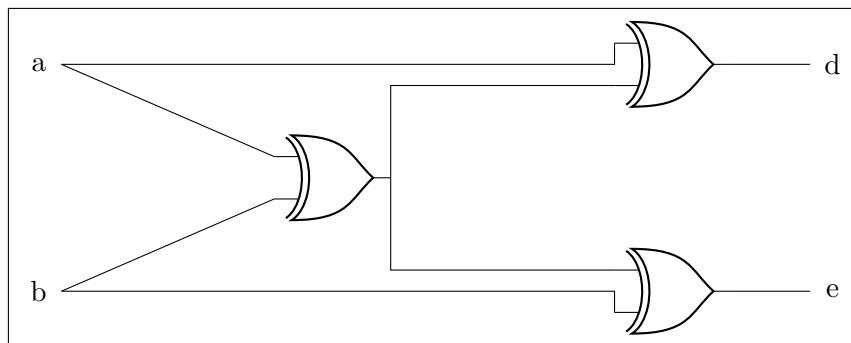
On peut maintenant montrer le résultat demandé :

$$d = a \oplus c = a \oplus (a \oplus b) = a \oplus a \oplus b = b$$

$$e = b \oplus c = b \oplus (a \oplus b) = b \oplus a \oplus b = a$$

■

b. On donne ci-dessous le circuit correspondant :



1.4 Exercice 4

a. On écrit le code suivant :

```

1 int main()
2 {
3     printf("Sizeof int: %lu octets\n", sizeof(int));
4     printf("Sizeof short: %lu octets\n", sizeof(short));
5     printf("Sizeof char: %lu octets\n", sizeof(char));
6     return 0;
7 }
  
```

La sortie est la suivante :

```

1 Sizeof int: 4 octets
2 Sizeof short: 2 octets
3 Sizeof char: 1 octets
  
```

b. On écrit le code suivant :

```

1 int main()
2 {
3     int a = pow(2, 31);
4     int b = pow(2, 31);
5     int c = a + b;
6     printf("%d\n", c);
7     return 0;
8 }
  
```

La sortie affiche 0, ce qui correspond bien à $2^{32} \bmod (2^{32})$

1.5 Exercice 5

On donne ci-dessous l'écriture binaire sur 4 et 8 bits de 0, 1, -1 et -2 :

x	4 bits	8 bits
0 :	0000	0000 0000
1 :	0001	0000 0001
-1 :	1111	1111 1111
-2 :	1110	1111 1110

1.6 Exercice 6

- $m_1 = 0001$ et $m_{-1} = 1001$.
- En abusant de la notation + pour des mots : $m_0 = m_1 + m_{-1} = 1010$.
- En suivant la règle de signes, 1010 est l'encodage de -2.

1.7 Exercice 7

Soit b un nombre de bits. Soit x un entier relatif qu'on souhaite représenter sur b bits.

Si $x \geq 0$, alors l'encodage de x correspond à une écriture dans $[0, 2^{b-1} - 1]$, alors cette écriture commence par un zéro (de 00...0 à 01...1). Si $x < 0$, alors $2^b - x \in [2^{b-1}, 2^b - 1]$ (soit de 10...0 à 11...1), son écriture commence par un 1.

■

1.8 Exercice 8

Dans le premier code, on dispose de 2 cases mémoires différentes. Le résultat affiché est -106 pour la valeur de d , ce qui est normal puisque d est signé.

Dans le deuxième code, on utilise une seule case mémoire à travers l'utilisation de deux pointeurs, un signé et un non signé. Le résultat affiché est identique au premier code.

Cela permet de montrer que la mémoire est "non typée", l'interprétation de la valeur mémoire dépend directement du type de l'objet qui lit cette valeur.

1.9 Exercice 9

- 10 s'écrit 2×5 et toute puissance de 2 s'écrit 2^k où $k \in \mathbb{N}$. Ainsi, si $x \in 2^{\mathbb{N}}$ (par abus de langage) est divisible par 10, alors x contient au moins 2 et 5 dans sa décomposition en facteurs premiers, ce qui donne une contradiction avec la propriété précédemment énoncée.

■

b. Supposons que 0.1 soit représentable sur kl bits. Alors, d'après le résultat du cours, $2^l \times 0.1$ est un entier, autrement dit 2^l est divisible par 10. D'après la question précédente, c'est impossible. ■

1.10 Exercice 10

L'écriture binaire approchée de 0.1 est $0.0001\ 1001_2$, de valeur décimale 0.09765625.

1.11 Exercice 11

a. On écrit la fonction suivante en C :

```
1 int main() {
2     for (int i = 0; i < 10; i++){
3         for (int j = 0; j < 10; j++){
4             float a = i/10.0, b = j/10.0, c = (i+j)/10.0;
5             printf("(%d, %d) : %s\n", i, j, (a+b == c)?"true":
6                 "false");
7         }
8     }
9     return 0;
}
```

La sortie affichée contient, entre autres, les résultats suivants :

- (1, 4) : true
- (1, 5) : true
- (1, 6) : false
- (1, 7) : true
- (1, 8) : false
- (3, 4) : false
- (3, 5) : true
- (3, 6) : false
- (3, 7) : true
- (3, 8) : true

b. On modifie seulement la ligne d'affichage dans le code ¹ :

```
1 printf("%.16f + %.16f = %.16f\n", a, b, c);
```

On donne seulement les résultats pour les 5 premiers couples ci-dessus :

```
1 0.1000000014901161 + 0.4000000059604645 = 0.5000000000000000
2 0.1000000014901161 + 0.5000000000000000 = 0.6000000238418579
3 0.1000000014901161 + 0.6000000238418579 = 0.6999999880790710
4 0.1000000014901161 + 0.6999999880790710 = 0.8000000119209290
5 0.1000000014901161 + 0.8000000119209290 = 0.8999999761581421
```

c. On remarque que pour une addition, l'égalité $a+b==c$ est vérifiée lorsque a et b sont représentables, ou quand l'un des deux seulement l'est. Dans le second cas, l'erreur de représentation n'a pas eu d'impact sur le résultat puisque c'est la seule erreur du calcul. Ainsi l'égalité reste vraie.

1. Comme un `int` est codé sur 4 octets, on donne 16 caractères à chaque affichage.

En revanche, dès que les deux flottants ne sont pas représentables, les erreurs s'accumulent et alors la représentation de `c` peut différer de la valeur de `a+b`.

Dans un cas général si on prend $x, y \in \mathbb{R}$ tels que $x = y$, on aura `x==y` quand x et y sont représentables sur un nombre de bits donné². Dans les autres cas, il est possible d'obtenir un résultat correct mais cela résulte plutôt du hasard.

d. On modifie la fonction précédente :

```

1 int main() {
2     for (int i = 0; i < 10; i++){
3         for (int j = 0; j < 10; j++){
4             float a = i/10.0, b = j/10.0, c = (i+j)/10.0;
5             if (a+b!=c){
6                 int m1 = a+b>c;
7                 int m2 = a+b+.000001>c+.000001;
8                 printf("(%d, %d) : %d %d\n", i, j, m1, m2);
9             }
10        }
11    }
12    return 0;
13 }
```

On obtient la sortie suivante :

```

1 // a+b > c is tested before and after adding 1e-6
2 (1, 6) : 1 1
3 (1, 8) : 1 1
4 (3, 4) : 1 1
5 (3, 6) : 1 1
6 (4, 3) : 1 1
7 (6, 1) : 1 1
8 (6, 3) : 1 1
9 (6, 8) : 1 1
10 (7, 9) : 0 0
11 (8, 1) : 1 1
12 (8, 6) : 1 1
13 (9, 7) : 0 0
```

On constate que l'ordre est conservé à chaque fois. Comme 10^{-6} n'est pas représentable sur 16 bits ($10^{-6} \approx 2^{-20}$), on simule l'effet de la propagation d'une erreur.

On en déduit que des erreurs successives d'arrondi ne bousculent pas l'ordre sur des valeurs arrondies. Toutefois ici on effectue le même calcul des deux côtés. On peut donc toujours les comparer³, même après plusieurs calculs, puisque l'ordre est conservé. On peut également penser que cela ne provoquera pas d'évolution chaotique ou aléatoire de ces valeurs dans les calculs.

2. Il faut tout de même faire attention à la précision. Augmenter la précision ne rendra pas les calculs exacts pour autant.

3. Pas l'égalité.

En revanche, pour une suite d'opérations inconnue, cela risque de devenir insignifiant de vouloir comparer deux flottants.

1.12 Exercice 12

On note $s_1 = \sum_{i=1}^k \frac{1}{i}$ et $s_2 = \sum_{i=k}^1 \frac{1}{i}$.

On écrit la fonction suivante :

```

1 int main() {
2     float s1 = 0.0, s2 = 0.0;
3     long k = 1000000000;
4     for (double i = 1; i < k+1; i++){
5         s1 += 1/i;
6         s2 += 1/(k+1-i);
7     }
8     printf("k=%ld\n    s1 = %.16f\n    s2 = %.16f\n", k, s1, s2)
9     ;
10    return 0;
11 }
```

Les résultats sont les suivants :

```

1 k=1000
2     s1 = 7.4854784011840820
3     s2 = 7.4854717254638672
4 k=1000000
5     s1 = 14.3573579788208008
6     s2 = 14.3926515579223633
7 k=1000000000
8     s1 = 15.4036827087402344
9     s2 = 18.8079185485839844
```

On constate que les résultats diffèrent⁴ d'autant plus que le nombre d'erreurs successives est grand, ce qui n'est pas étonnant. Là où le résultat interpelle, c'est que l'ordre de parcours de la somme a un impact conséquent sur le résultat.

Cette erreur est due au fait que lorsque le parcours est décroissant, on finit pas les petits nombres. Ces derniers causent alors des erreurs d'arrondi car leur ordre de grandeur est faible par rapport aux premiers nombres.

1.13 Exercice 13

- a. La taille d'un `float` en C est de 4 *bytes*, soit 32 bits :
 - 1 bit de signe
 - 8 bits d'exposant

4. Même plus que ça, on dirait que la série converge! On connaît l'équivalent pour la série harmonique : $H_n \underset{n \rightarrow +\infty}{\sim} \gamma + \ln(n)$ où $\gamma \approx 0.5$ est la constante d'Euler. Pour $n = 10^9$, on devrait donc plutôt être autour de $H_n \approx 21$.

- 23 bits de mantisse
- b. Convertissons d'abord $0x414BD000$ en binaire :

$$414BD000_{16} = 01000001010010111101000000000000_2$$

On identifie ensuite les bits de signe, d'exposant et de mantisse :

$$\underbrace{0}_S \underbrace{10000010}_{E=130} \underbrace{100101111010000000000000}_T$$

Le biais b valant 127, notre exposant ici vaut $E - b = 3$. Donc, en "écriture binaire à virgule", on obtient :

$$1.10010111101 \times 2^3 = \underbrace{1100}_{=12} . \underbrace{10111101}_{=0.73828125}$$

Soit finalement :

$$0x414BD000_{16} \text{ encode } 12.73828125$$

c. En suivant le raisonnement inverse, on peut trouver l'exposant et la mantisse de l'encodage de 0.1. On commence par l'écrire en "binaire à virgule" :

$$0.1 = 000110011001100110011001101_2$$

On décale la virgule pour trouver l'exposant :

$$0.1 = 1.10011001100110011001101_2 \times 2^{-4}$$

Cela donne donc un exposant de $E = b - 4 = 123 = 01111011_2$ sur 8 bits. Enfin, $0.1 > 0$ donc on met un premier bit à 0. On obtient donc l'encodage suivant – dont on donne également la valeur décimale réelle – pour le nombre 0.1 :

$$\underbrace{0}_{\text{signe}} \underbrace{01111011}_{\text{exposant}} \underbrace{1001100110011001101}_{\text{mantisse}}_2 = 0.100000001490116119384765625$$

2 Chapitre 2

2.1 Exercice 1

a. Lorsqu'on crée par exemple un tableau de taille un milliard, on obtient une erreur similaire à la suivante :

```
1 [1] 54133 segmentation fault ./a.out
```

b. En expérimentant à la main, on trouve qu'on peut créer un tableau de taille maximale 2 096 286.

2.2 Exercice 2

On vérifie par exemple qu'on peut créer un tableau de taille un milliard.

2.3 Exercice 3

La machine utilisée pour ce TD utilise la convention *little endian*.

2.4 Exercice 4

On écrit déjà le code suivant :

```
1 #include <stdio.h>
2 #include <x86intrin.h>
3
4 unsigned long int squareSum(int n){
5     unsigned long int tic, toc;
6     unsigned int ui;
7     int a = 0;
8     tic = __rdtscp(&ui);
9     for (int i=0; i < n; ++i){
10         a = a*a+a*a;
11     }
12     toc = __rdtscp(&ui);
13     return toc-tic;
14 }
15
16 int main(){
17     printf("n=%d: %lu tics\n", 1000, squareSum(1000));
18     printf("n=%d: %lu tics\n", 10000, squareSum(10000));
19     printf("n=%d: %lu tics\n", 1000000, squareSum(1000000));
20     printf("n=%d: %lu tics\n", 10000000, squareSum(10000000));
21     printf("n=%d: %lu tics\n", 100000000, squareSum(100000000));
22     return 0;
23 }
```

On obtient les résultats suivants :

n	10^3	10^4	10^6	10^7	10^8
mesure	3310	42456	7134866	67719612	637584056

On note qu'on semble bien obtenir une relation qui a l'air linéaire, hormis la dernière mesure.

2.5 Exercice 5

On obtient les résultats suivants :

appel\ n	10^3	10^5	10^7	10^8	10^9
1er	9108	835910	52244848	532734614	5350923602
2ème	3238	685744	31633110	317468584	3161147144
3ème	3088	322012	32683670	341747932	3100426690
$\frac{\text{écart max.}}{\text{moyenne}}$ en %	117	83	53	53	58

2.6 Exercice 6

a. On teste la fonction avec la fonction `test` calculant la somme des premiers carrés. On obtient des résultats cohérents (croissance linéaire). La fonction `print_timing` comprend une amorce qui exécute 10 fois la fonction à tester et moyenne les mesures sur 100 appels.

b. Voici la fonction `print_timing` :

```
1 void print_timing(int arg, void (*func)(int), int nb_boot, int
  nb_call)
2 {
3     // boot up
4     for (int i = 0; i < nb_boot; i++)
5     {
6         func(arg);
7     }
8
9     // measure
10    unsigned long int tic, toc;
11    unsigned int ui;
12    tic = _rdtscp(&ui);
13    for (int i = 0; i < nb_call; i++)
14    {
15        func(arg);
16    }
17    toc = _rdtscp(&ui);
18
19    printf("average time : %lu\n", (toc - tic) / n);
20 }
```

c. Oui.

d. On obtient les résultats suivants :

appel\ <i>n</i>	10^3	10^5	10^7	10^8	10^9
mesure <code>print_timing</code>	4214	432308	31686176	308206789	3176004203

2.7 Exercice 7

On obtient les résultats suivants (arrondis) :

pas	1	2	3	4	8	16	32
mesure (10^6 tics)	59	61	63	66	83	102	197

2.8 Exercice 8