

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

# Informatique

## 10

# Représentation des nombres

### *Résumé*

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

## Représentation des nombres

Binaires entiers																																																																																																																																
Définition	Un nombre codé en binaire sur $n$ bits est écrit sous la forme d'une chaîne de 0 et 1. On appelle mot $m$ l'ensemble des bits représentant le nombre en binaire : $m = 000110110$ On trouve parfois l'écriture $m = 0b000110110$ pour spécifier que c'est un binaire																																																																																																																															
Propriétés	Un entier codé sur $n$ bits peut prendre $2^n$ valeurs différentes (2 possibilités par bit), 0 inclus Le plus grand entier représenté sur $n$ bits vaut donc $2^n - 1$ Exemple sur 8 bits : $255_{(10)} = 11111111_{(2)} = 2^8 - 1$																																																																																																																															
Les premiers entiers		Base 10	Base 2	Base 10	Base 2	Base 10	Base 2																																																																																																																									
		0	0	6	110	11	1100																																																																																																																									
		1	1	7	111	13	1101																																																																																																																									
		2	10	8	1000	14	1110																																																																																																																									
		3	11	9	1001	15	1111																																																																																																																									
		4	100	10	1010	16	10000																																																																																																																									
		5	101	11	1011	17	10001																																																																																																																									
Transcodage	Base 10 → Binaire	Diviser par 2 l'entier et diviser le quotient obtenu jusqu'à ce qu'il soit nul. Remonter les restes de chaque division pour créer le mot binaire.																																																																																																																														
		<table><tr><td></td><td>1000</td><td>2</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td><math>n_0</math></td><td>0</td><td>500</td><td>2</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td><math>n_1</math></td><td>0</td><td>250</td><td>2</td><td></td><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td><math>n_2</math></td><td>0</td><td>125</td><td>2</td><td></td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td><math>n_3</math></td><td>1</td><td>62</td><td>2</td><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td><math>n_4</math></td><td>0</td><td>31</td><td>2</td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td><td><math>n_5</math></td><td>1</td><td>15</td><td>2</td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td><td></td><td><math>n_6</math></td><td>1</td><td>7</td><td>2</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td><math>n_7</math></td><td>1</td><td>3</td><td>2</td></tr><tr><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td><math>n_8</math></td><td>1</td><td>1</td></tr><tr><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td><math>n_9</math></td><td>1</td></tr></table> <div>&gt;&gt;&gt; 0b111101000      1000<sub>(10)</sub> = 111101000<sub>(2)</sub>      &gt;&gt;&gt; bin(1000) 1000      '0b111101000'</div>								1000	2									$n_0$	0	500	2									$n_1$	0	250	2									$n_2$	0	125	2									$n_3$	1	62	2									$n_4$	0	31	2									$n_5$	1	15	2									$n_6$	1	7	2									$n_7$	1	3	2									$n_8$	1	1										$n_9$
		1000	2																																																																																																																													
$n_0$		0	500	2																																																																																																																												
		$n_1$	0	250	2																																																																																																																											
			$n_2$	0	125	2																																																																																																																										
				$n_3$	1	62	2																																																																																																																									
					$n_4$	0	31	2																																																																																																																								
						$n_5$	1	15	2																																																																																																																							
							$n_6$	1	7	2																																																																																																																						
								$n_7$	1	3	2																																																																																																																					
								$n_8$	1	1																																																																																																																						
									$n_9$	1																																																																																																																						
Binaire → Base 10	Il suffit de sommer les puissances de 2 associées à chaque bit $111101000_{(2)}$ $= 1 \cdot 2^9 + 1 \cdot 2^8 + 1 \cdot 2^7 + 1 \cdot 2^6 + 1 \cdot 2^5 + 0 \cdot 2^4 + 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0$ $= 2^9 + 2^8 + 2^7 + 2^6 + 2^5 + 2^3 = 512 + 256 + 128 + 64 + 32 + 8$ $111101000_{(2)} = 1000_{(10)}$																																																																																																																															
Entiers multi-précision		Lorsque les entiers dépassent la plus grande valeur stockable de $2^n - 1$ , ils sont codés en multi-précision, ce qui revient à les représenter par une liste de leurs termes ( $N = 111; L = [1,1,1]$ ), puis d'utiliser des algorithmes spécifiques réalisant les opérations arithmétiques de base (addition, soustraction, multiplication et division). La seule limite est alors la taille mémoire utilisée pour les stocker.																																																																																																																														

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

Binaires réels																														
Principe	En base 10, lorsque nous manipulons des nombres à virgule : <ul style="list-style-type: none"><li>- Les chiffres avant la virgule sont des puissances de 10</li><li>- Les chiffres après la virgule sont des puissances de 1/10</li></ul> C'est la même chose en base 2, en puissances de 2 ou 1/2																													
Transcodage	Base 10 → Binaire	<p>Décomposer le réel en sa partie entière et sa partie décimale : <math>5,375 = 5 + 0,375</math></p> <p>Transcoder la partie entière comme vu précédemment : <math>5_{(10)} = 101_{(2)}</math></p> <p>Transcoder la partie décimale de la sorte :</p> <ul style="list-style-type: none"><li>- La multiplier par 2. Séparer l'entier avant la virgule (0 ou 1) de la nouvelle partie décimale</li><li>- Recommencer jusqu'à obtenir une partie décimale nulle (si existe)</li><li>- Récupérer la partie décimale binaire</li></ul> <table><thead><tr><th colspan="2">Principe</th><th colspan="2">Présentation améliorée</th></tr></thead><tbody><tr><td><math>0,375 * 2 = 0,750</math> <i>Bit 0</i> → { <i>Nouvelle partie décimale 0,75</i></td><td><math>0,750 * 2 = 1,5</math> <i>Bit 1</i> → { <i>Nouvelle partie décimale 0,5</i></td><td><table><tr><td>0,375</td><td>* 2 =</td><td>0,</td><td>750</td></tr><tr><td>0,750</td><td>* 2 =</td><td>1,</td><td>5</td></tr><tr><td>0,50</td><td>* 2 =</td><td>1,</td><td>0</td></tr></table></td><td></td></tr><tr><td colspan="2"><math>0,50 * 2 = 1</math> <i>Bit 1</i> → { <i>Nouvelle partie décimale 0</i></td><td colspan="2"></td></tr><tr><td colspan="4"><math>(0,375)_{10} = (0,011)_2</math></td></tr></tbody></table> <p>Finalement : <math>(5,375)_{10} = (101,011)_2</math></p> <p>Remarque : on peut réaliser ce transcodage en notation scientifique binaire : <math>(0,011)_2 = (11)_2 \cdot 2^{-3} = \frac{3}{8} = 0,375</math></p>	Principe		Présentation améliorée		$0,375 * 2 = 0,750$ <i>Bit 0</i> → { <i>Nouvelle partie décimale 0,75</i>	$0,750 * 2 = 1,5$ <i>Bit 1</i> → { <i>Nouvelle partie décimale 0,5</i>	<table><tr><td>0,375</td><td>* 2 =</td><td>0,</td><td>750</td></tr><tr><td>0,750</td><td>* 2 =</td><td>1,</td><td>5</td></tr><tr><td>0,50</td><td>* 2 =</td><td>1,</td><td>0</td></tr></table>	0,375	* 2 =	0,	750	0,750	* 2 =	1,	5	0,50	* 2 =	1,	0		$0,50 * 2 = 1$ <i>Bit 1</i> → { <i>Nouvelle partie décimale 0</i>				$(0,375)_{10} = (0,011)_2$			
	Principe		Présentation améliorée																											
$0,375 * 2 = 0,750$ <i>Bit 0</i> → { <i>Nouvelle partie décimale 0,75</i>	$0,750 * 2 = 1,5$ <i>Bit 1</i> → { <i>Nouvelle partie décimale 0,5</i>	<table><tr><td>0,375</td><td>* 2 =</td><td>0,</td><td>750</td></tr><tr><td>0,750</td><td>* 2 =</td><td>1,</td><td>5</td></tr><tr><td>0,50</td><td>* 2 =</td><td>1,</td><td>0</td></tr></table>	0,375	* 2 =	0,	750	0,750	* 2 =	1,	5	0,50	* 2 =	1,	0																
0,375	* 2 =	0,	750																											
0,750	* 2 =	1,	5																											
0,50	* 2 =	1,	0																											
$0,50 * 2 = 1$ <i>Bit 1</i> → { <i>Nouvelle partie décimale 0</i>																														
$(0,375)_{10} = (0,011)_2$																														
Binaire → Base 10	<p>Il suffit de sommer les puissances de 2 associées à chaque bit</p> $101,011_{(2)}$ $= 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3}$ $= 4 + 1 + \frac{1}{4} + \frac{1}{8} = 5 + 0,25 + 0,125 = 5 + 0,375 = 5,375$																													

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

## Entiers relatifs

### Principe

Associer à  $2^n$  entiers naturels  $N$  codés de manière classique sur  $n$  bits  $2^n$  entiers relatifs  $Z$  selon une formule  $N = f(Z)$  ou  $N = f^{-1}(Z)$ .

### Complément à 2

Le principe consiste à utiliser un entier naturel  $N$  (pour  $n = 3$  bits,  $N \in [0,7]$ ), et de lui associer un entier relatif  $Z$  tel que :

$$N = \begin{cases} Z & \text{si } Z \geq 0 \\ 2^n - |Z| = 2^n + Z & \text{si } Z < 0 \end{cases} \Leftrightarrow Z = \begin{cases} N & \text{si } N \leq 2^{n-1} - 1 \\ N - 2^n & \text{si } N > 2^{n-1} - 1 \end{cases}$$

$$N \in [0, 2^n] ; \quad Z \in [-2^{n-1}, 2^{n-1} - 1]$$

$(Z)_{10}$	-4	-3	-2	-1	0	1	2	3
$(N)_{10}$	4	5	6	7	0	1	2	3
$(N)_2$	$(100)_2$	$(101)_2$	$(110)_2$	$(111)_2$	$(000)_2$	$(001)_2$	$(010)_2$	$(011)_2$

Remarque : On voit que le premier bit est un bit de signe

Avantage : La somme binaire de deux entiers relatif  $Z$  fonctionne 😊

$$(-3)_{10} + (1)_{10} \Leftrightarrow (101)_2 + (001)_2 = (110)_2 \Leftrightarrow (-2)_{10}$$

Exemple de transcoding sur 8 bits :

$(N)_2$	$(N)_{10}$	$(Z)_{10}$
$(11010111)_2$	$(215)_{10}$	$(-41)_{10}$

### Codage par excès

Le principe consiste à utiliser un entier naturel  $N$  (pour  $n = 3$  bits,  $N \in [0,7]$ ), et de lui associer un entier relatif  $Z$  « décalé ». Deux décalages semblent logiques :

- Décaler de 3 : 0 devient -3 et 7 devient 4
- Décaler de 4 : 0 devient -4 et 7 devient 3

Le choix est fait de procéder au premier décalage dont l'expression en fonction de  $n$  s'écrit :

$$B = 2^{n-1} - 1$$

$$Z = N - B \Leftrightarrow N = Z + B$$

$$N \in [0, 2^n] ; \quad Z \in [-2^{n-1} + 1, 2^{n-1}]$$

$(Z)_{10}$	-3	-2	-1	0	1	2	3	4
$(N)_{10}$	0	1	2	3	4	5	6	7
$(N)_2$	$(000)_2$	$(001)_2$	$(010)_2$	$(011)_2$	$(100)_2$	$(101)_2$	$(110)_2$	$(111)_2$

Exemple de transcoding sur 8 bits ( $B = 127$ ) que nous utiliserons plus tard :

$N$	$Z$
0	-127
255	128

$(N)_2$	$(N)_{10}$	$(Z)_{10}$
$(10101100)_2$	$(86)_{10}$	$(-41)_{10}$

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

Virgule flottante – Norme IEEE 754																											
Contexte	Il est impossible de représenter une infinité de nombres dans l’ordinateur. Une représentation officielle a été choisie, permettant de représenter les réels en machine à l’aide de « flottants », ensemble fini de nombres décimaux																										
Exemple	<table><tr><td colspan="4">–289456,5 en base 10</td></tr><tr><td colspan="4">–2,894565. 10<sup>5</sup></td></tr><tr><td>–</td><td>2</td><td>894565</td><td>5</td></tr><tr><td>Signe</td><td>Caractéristique</td><td>Mantisse</td><td>Puissance</td></tr><tr><td colspan="4">Partie significative</td></tr></table>							–289456,5 en base 10				–2,894565. 10 <sup>5</sup>				–	2	894565	5	Signe	Caractéristique	Mantisse	Puissance	Partie significative			
	–289456,5 en base 10																										
	–2,894565. 10 <sup>5</sup>																										
	–	2	894565	5																							
	Signe	Caractéristique	Mantisse	Puissance																							
	Partie significative																										
	<table><tr><td colspan="4">–289456,5 en virgule flottante</td></tr><tr><td colspan="4">–(1,0001101010101100001)<sub>2</sub>. 2<sup>18</sup></td></tr><tr><td>–</td><td>1</td><td>0001101010101100001</td><td>18</td></tr><tr><td>Signe</td><td>Caractéristique</td><td>Mantisse</td><td>Puissance</td></tr><tr><td colspan="4">Partie significative</td></tr></table>							–289456,5 en virgule flottante				–(1,0001101010101100001) <sub>2</sub> . 2 <sup>18</sup>				–	1	0001101010101100001	18	Signe	Caractéristique	Mantisse	Puissance	Partie significative			
	–289456,5 en virgule flottante																										
	–(1,0001101010101100001) <sub>2</sub> . 2 <sup>18</sup>																										
	–	1	0001101010101100001	18																							
Signe	Caractéristique	Mantisse	Puissance																								
Partie significative																											
La caractéristique est toujours égale à 1 sauf le cas exceptionnel de 0																											
Le bit de signe vaut $\begin{cases} s = 0 \Leftrightarrow x > 0 \\ s = 1 \Leftrightarrow x < 0 \end{cases}$																											
La puissance (entier relatif Z), est codée sur n bits (8 en simple précision) en un entier naturel N : $N = Z + (2^{n-1} - 1)$ . Dans ce cas, $18 + 127 = (145)_{10} = (10010001)_2$																											
<table><tr><td>1</td><td>10010001</td><td>0001101010101100001</td></tr><tr><td>Signe</td><td>Puissance</td><td>Mantisse</td></tr></table>							1	10010001	0001101010101100001	Signe	Puissance	Mantisse															
1	10010001	0001101010101100001																									
Signe	Puissance	Mantisse																									
En simple précision, la mantisse est codée sur 23 bits, on complète donc de quelques 0 :																											
<table><tr><td rowspan="2">Représentation</td><td>1</td><td>10010001</td><td>00011010101011000010000</td></tr><tr><td>Signe</td><td>Puissance</td><td>Mantisse</td></tr><tr><td>32 bits</td><td>1 bit</td><td>8 bits</td><td>23 bits</td></tr></table>							Représentation	1	10010001	00011010101011000010000	Signe	Puissance	Mantisse	32 bits	1 bit	8 bits	23 bits										
Représentation	1	10010001	00011010101011000010000																								
	Signe	Puissance	Mantisse																								
32 bits	1 bit	8 bits	23 bits																								
Le nombre –289456,5 est codé en virgule flottante simple précision (32 bits): 11001000100011010101011000010000																											
Norme																											
		Nombre de bits	Bit de signe	Bits exposant	Bits implicite	Bits mantisse	Décalage																				
	Simple précision	32	1	8	1	23	127																				
	Double précision	64	1	11	1	52	1023																				
	Quadruple précision	128	1	15	1	112	16383																				
Remarques	Certains nombres sont réservés, il existe une représentation normalisée et une représentation dénormalisée (cf cours), etc.. Nous n’abordons ici que la structure des nombres pour comprendre qu’ils puissent avoir des conséquences sur nos calculs.																										

Dernière mise à jour	Informatique	Denis DEFAUCHY – <a href="#">Site web</a>
23/05/2022	10 - Représentation des nombres	Résumé

Transcodage En 32 bits	Virgule flottante 32 bits → Réel	<div>11000001001001110000000000000000</div> <table><tr><th colspan="3">Identification des 3 parties</th></tr><tr><td>1</td><td>10000010</td><td>010011100000000000000000</td></tr><tr><td>Signe</td><td>Puissance</td><td>Mantisse</td></tr></table> <p>Le signe est négatif ; <math>(10000010)_2 = (130)_{10}</math> <math>N = 130</math> ; <math>Z = N - (2^{n-1} - 1) = 130 - 127 = 3</math> Sans oublier d'ajouter le bit implicite, la partie significative vaut : <math>(1,0100111)_2 = 1,3046875</math> ; <math>x = -1,3046875 \cdot 2^3 = -10,4375</math></p>				Identification des 3 parties			1	10000010	010011100000000000000000	Signe	Puissance	Mantisse	
	Identification des 3 parties														
1	10000010	010011100000000000000000													
Signe	Puissance	Mantisse													
	Réel → Virgule flottante 32 bits	<p><math>x = (-10,125)_{10} = (-1010,001)_2 = (-1,010001)_2 \cdot 2^3</math> <math>Z = 127</math> ; <math>N = 3 + (2^{n-1} - 1) = 3 + 127 = 130</math> <math>(130)_{10} = (10000010)_2</math></p> <table><tr><th colspan="3"><math>(-10,125)_{10}</math></th></tr><tr><td>1</td><td>10000010</td><td>010001000000000000000000</td></tr><tr><td>Signe</td><td>Puissance</td><td>Mantisse</td></tr></table> <div>11000001001000100000000000000000</div>				$(-10,125)_{10}$			1	10000010	010001000000000000000000	Signe	Puissance	Mantisse	
$(-10,125)_{10}$															
1	10000010	010001000000000000000000													
Signe	Puissance	Mantisse													
Limites (cf cours)	Justesse	C'est la correspondance entre nombre représenté et réel associé). On peut représenter une quantité finie de nombres : <table><tr><td>32 bits</td><td colspan="3">4 294 967 296</td></tr><tr><td>64 bits</td><td colspan="3">18 446 744 073 710 000 000</td></tr></table> <p>La justesse est assurée si le nombre de bits de la mantisse suffit à la représenter complètement.</p>				32 bits	4 294 967 296			64 bits	18 446 744 073 710 000 000				
		32 bits	4 294 967 296												
	64 bits	18 446 744 073 710 000 000													
	Min Max	Valeurs normalisées	Valeur min normalisée	Valeur min dénormalisée	Valeur max										
Simple précision		$1,2 \cdot 10^{-38}$	$1,4 \cdot 10^{-45}$	$3,4 \cdot 10^{38}$											
	Double précision	$2,2 \cdot 10^{-308}$	$4,9 \cdot 10^{-324}$	$1,8 \cdot 10^{308}$											
	Ecart entre nombres successifs et applications en virgule flottante	<b>ODG</b> de l'écart $\Delta V$ entre $A = a \cdot 10^n$ avec $a \in [0 ; 10[$ et le suivant : <table><tr><td>Simple précision</td><td>Double précision</td><td><math>7 \approx 23 \log 2</math></td></tr><tr><td><math>10^{n-7}</math></td><td><math>10^{n-16}</math></td><td><math>16 \approx 52 \log 2</math></td></tr></table> <p>Une erreur d'arrondi peut conduire à <math>A = B + \Delta V</math> alors que A et B devraient être égaux → Etudier <math> A - B  &lt; \varepsilon</math>, avec <math>\varepsilon</math> de l'ordre de grandeur de <math>\Delta V</math> du tableau précédent... Attention, c'est un <b>ODG</b> ! On prendra <math>\varepsilon = 10^{n-6}</math> (32) ou <math>\varepsilon = 10^{n-15}</math> (64) par exemple <math>\varepsilon</math> trop faible <math>\Rightarrow</math> Le test dit <math>A \neq B</math> alors que <math>A = B</math> <math>\varepsilon</math> trop grand <math>\Rightarrow</math> Le test dit <math>A = B</math> alors que <math>A \neq B</math></p> <b>ODG</b> du nombre de chiffres significatifs (caractéristique + mantisse) en base 10, auxquels on peut avoir confiance : <table><tr><td>Simple précision</td><td>Double précision</td></tr><tr><td>7</td><td>16</td></tr></table>				Simple précision	Double précision	$7 \approx 23 \log 2$	$10^{n-7}$	$10^{n-16}$	$16 \approx 52 \log 2$	Simple précision	Double précision	7	16
Simple précision		Double précision	$7 \approx 23 \log 2$												
$10^{n-7}$		$10^{n-16}$	$16 \approx 52 \log 2$												
Simple précision		Double précision													
7	16														
Conséquences	Les erreurs d'arrondis ou l'overflow (nombre non représentable, induisant le codage d'une mauvaise valeur) ont déjà eu de lourdes conséquences (explosion de la fusée Ariane 5, Missile Patriot passant à travers un système anti-missile...) Attention au test == – Attention aux sommes répétées 1 et 1.0 sont fondamentalement différents – Attention à l'évaluation d'une dérivée														
	Numpy	<pre>a = np.float32(0.1) - b = np.uint8(10) A = np.array([1,2], dtype='int8')</pre>													