

# BROUILLON - CANDIDAT - QUELS DÉCALAGES ENTRE LES ENTIERS SIGNÉS ET LES NON SIGNÉS ?

CHRISTOPHE BAL

*Document, avec son source L<sup>A</sup>T<sub>E</sub>X, disponible sur la page  
<https://github.com/bc-writing/drafts>.*

## Mentions « légales »

Ce document est mis à disposition selon les termes de la licence Creative Commons “Attribution - Pas d’utilisation commerciale - Partage dans les mêmes conditions 4.0 International”.



## TABLE DES MATIÈRES

1. Entiers non signés	1
2. Entiers signés	2
2.1. Comment ça marche ?	2
2.2. Qu’est-ce qui motive ces choix ?	2
2.3. Pourquoi ça marche ?	4
3. Division par 2 et décalage à droite	6
3.1. Entiers non signés	6
3.2. Entiers signés	6
3.3. Synthèse	6
4. Multiplication par 2 et décalage à gauche	7
4.1. Entiers non signés	7
4.2. Entiers signés	7
4.3. Synthèse	8

## 1. ENTIERS NON SIGNÉS

Dans ce document, nous allons travailler uniquement avec des entiers non signés, c’est à dire des naturels dont l’écriture en base 2 tient sur 8 bits même si les raisonnements tenus se généralisent à 16, 32 ou 64 bits. Ce choix tient juste à des contraintes de lisibilité du document.

Voici des exemples de représentations de type little-endian : on part à gauche de la puissance la plus haute, pour nous ce sera toujours  $2^7$ , puis on va vers la plus basse (*c’est le même principe que l’écriture des naturels en base 10*).

1	0	0	0	0	1	0	1
---	---	---	---	---	---	---	---

 $\leftarrow 133 = 128 + 4 + 1 = 2^7 + 2^2 + 2^0$

0	0	1	0	0	1	1	0
---	---	---	---	---	---	---	---

 $\leftarrow 38 = 32 + 4 + 2 = 2^5 + 2^2 + 2^1$

Avec 8 bits on peut représenter les naturels de 0 à  $1 + 2 + 2^2 + \dots + 2^7 = 2^8 - 1 = 255$ .

## 2. ENTIERS SIGNÉS

2.1. **Comment ça marche ?** Voici des exemples de représentations toujours sur 8 bits et de type little-endian où la case rouge indique le bit utilisé pour coder le signe.

$$\boxed{0} \boxed{1} \boxed{1} \boxed{1} \boxed{1} \boxed{0} \boxed{0} \boxed{0} \leftarrow 120 = 64 + 32 + 16 + 8 = \boxed{0 \times (-2^7)} + 2^6 + 2^5 + 2^4 + 2^3$$

$$\boxed{1} \boxed{0} \boxed{0} \boxed{0} \boxed{1} \boxed{0} \boxed{0} \boxed{0} \leftarrow -120 = -128 + 8 = \boxed{1 \times (-2^7)} + 2^3$$

$$\boxed{1} \boxed{0} \boxed{0} \boxed{0} \boxed{0} \boxed{0} \boxed{0} \boxed{0} \leftarrow -128 = \boxed{1 \times (-2^7)} + 0$$

Avec 8 bits on peut représenter les naturels de  $-2^7 = -128$  à  $1 + 2 + 2^2 + \dots + 2^6 = 2^7 - 1 = 127$ .

Ce choix permet d'effectuer des additions relatives toujours sous forme d'addition bit à bit, et donc par conséquent de faire aussi des soustractions via des additions bit à bit. Voici quelques exemples sans dépassement de capacité.

- Addition de deux positifs.

$$\begin{array}{r} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{0} \leftarrow 84 = 64 + 16 + 4 \\ + \boxed{0} \boxed{0} \boxed{0} \boxed{1} \boxed{0} \boxed{0} \boxed{0} \boxed{1} \leftarrow 17 = 16 + 1 \\ = \boxed{0} \boxed{1} \boxed{1} \boxed{0} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \leftarrow 101 = 64 + 32 + 4 + 1 \end{array}$$

- Addition de deux négatifs (*la retenue finale hors capacité est ignorée*).

$$\begin{array}{r} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \boxed{0} \boxed{0} \leftarrow -84 = \boxed{-128} + 44 = \boxed{-128} + 32 + 8 + 4 \\ + \boxed{1} \boxed{1} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \boxed{1} \boxed{1} \leftarrow -17 = \boxed{-128} + 111 = \boxed{-128} + 64 + 32 + 8 + 4 + 2 + 1 \\ = \boxed{1} \boxed{0} \boxed{0} \boxed{1} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \leftarrow -101 = \boxed{-128} + 27 = \boxed{-128} + 16 + 8 + 2 + 1 \end{array}$$

- Addition de deux entiers de signes différents.

$$\begin{array}{r} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{0} \leftarrow 84 = 64 + 16 + 4 \\ + \boxed{1} \boxed{0} \boxed{0} \boxed{1} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \leftarrow -101 = \boxed{-128} + 27 = 16 + 8 + 2 + 1 \\ = \boxed{1} \boxed{1} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \boxed{1} \boxed{1} \leftarrow -17 = \boxed{-128} + 111 = -128 + 64 + 32 + 8 + 4 + 2 + 1 \end{array}$$

- Autre addition de deux entiers de signes différents (*la retenue finale hors capacité est ignorée*).

$$\begin{array}{r} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{0} \boxed{0} \leftarrow 84 = 64 + 16 + 4 \\ + \boxed{1} \boxed{1} \boxed{0} \boxed{0} \boxed{0} \boxed{0} \boxed{1} \boxed{1} \leftarrow -61 = \boxed{-128} + 67 = 64 + 2 + 1 \\ = \boxed{0} \boxed{0} \boxed{0} \boxed{1} \boxed{0} \boxed{1} \boxed{1} \boxed{1} \leftarrow 23 = 16 + 4 + 2 + 1 \end{array}$$

2.2. **Qu'est-ce qui motive ces choix ?** Pour comprendre comment découvrir ce type de procédé, nous allons raisonner en base 10 et imaginer que nous voulions calculer  $189 - 32 = 157$  en utilisant uniquement des additions. Tout un chacun sait effectuer ce calcul comme suit.

$$\begin{array}{r} 1 \ 8 \ 9 \\ - \quad 3 \ 2 \\ \hline = 1 \ 5 \ 7 \end{array}$$

Nous décidons d'ajouter 1000 à  $189 - 32$ . Nous obtenons :

$$\begin{array}{r} 1 \ 8 \ 9 \\ - \quad 3 \ 2 \\ + 1 \ 0 \ 0 \ 0 \\ \hline = 1 \ 1 \ 5 \ 7 \end{array}$$

Comme le résultat initial tenait sur les trois chiffres de droite, l'ajout de 1000 donne un nouveau résultat dont nous savons que les trois chiffres de droite donnent le résultat de  $189 - 32$  car la soustraction  $1157 - 1000$  revient à ignorer le 1 tout à droite.

D'autre part,  $1000 + 189 - 32 = 189 + 1000 - 32 = 189 + 968$  est une simple addition.

Notons enfin que  $1000 - 32 = 1 + 999 - 032 = 1 + 967$  où 9, 6 et 7 sont les compléments à 9 de 0, 3 et 2 respectivement.

En résumé, nous avons fait comme suit où • indique une ligne où a été fait le calcul d'un complément à 9 plus 1 (*nous avons encadré le signe final pour le différencier des signes d'opération*).

		1	8	9
−		0	3	2
<hr/>				
		1	8	9
+	●	9	6	8
=	1	1	5	7
<hr/>				
	<div>+</div>	1	5	7

Très bien mais que se passe-t-il avec  $32 - 189 = -157$  qui est négatif? Testons pour voir.

$$\begin{array}{r} \phantom{+} \phantom{\bullet} \phantom{8} \phantom{1} \phantom{1} \\ - \phantom{\bullet} \phantom{8} \phantom{1} \phantom{1} \phantom{1} \\ \hline \phantom{+} \phantom{\bullet} \phantom{8} \phantom{1} \phantom{1} \\ + \phantom{\bullet} \phantom{8} \phantom{1} \phantom{1} \phantom{1} \\ \hline = \phantom{\bullet} \phantom{8} \phantom{1} \phantom{1} \phantom{1} \end{array}$$

Oups ! Nous ne voyons pas directement 157. Normal ! Il reste à effectuer  $843 - 1000 = -157$ . Dans un tel cas, on note que  $843 - 1000 = -(1000 - 843)$  s'obtient via un complément à 9 plus 1 de 843 tout en ajoutant un signe moins. En résumé, nous procédons comme suit en notant qu'ici nous avons un 0 tout à gauche juste avant le résultat final contrairement au 1<sup>er</sup> cas où il y avait un 1.

			3	2
−		1	8	9
<hr/>				
			3	2
+	●	8	1	1
=	0	8	4	3
<hr/>				
	□ −	1	5	7

Ne nous emballons pas trop vite car il reste un cas problématique que nous n'avons pas abordé. Examinons ce qu'il se passe avec  $-32 - 189 = -221$ .

	—	0	3	2
—		1	8	9
<hr style="border: none; border-top: 1px solid black; margin: 2px 0;"/>				
	●	9	6	8
+	●	8	1	1
<hr style="border: none; border-top: 1px solid black; margin: 2px 0;"/>				
=		1	7	7

Nous devons affiner notre règle via la valeur du chiffre le plus à gauche car sinon ici nous aurions un résultat positif! En fait, il est facile de voir qu'il suffit de tenir compte du nombre de compléments à 9 plus 1 initiaux effectués<sup>1</sup> et du chiffre tout à gauche qui sera 0 ou 1. En effet, en notant  $n$  le nombre de compléments à 9 plus 1 effectués avant l'étape finale et  $g$  le chiffre à gauche de l'avant dernier résultat, nous avons les faits suivants.

1. C'est à dire le nombre de 1000 ajoutés.

- (1) **Cas 1** : pour  $189 - 32$ ,  $n = 1$ ,  $g = 1$  et nous avons juste lu le résultat en ignorant le chiffre tout à gauche sans faire de complément à 9 plus 1 supplémentaire.
- (2) **Cas 2** : pour  $32 - 189$ ,  $n = 1$ ,  $g = 0$  et nous avons fait un complément à 9 plus 1 supplémentaire, lu le résultat en ignorant le chiffre tout à gauche puis enfin ajouter un signe moins.
- (3) **Cas 3 actuel** : pour  $-32 - 189$ ,  $n = 2$ ,  $g = 1$  et nous devons faire comme dans le 2<sup>e</sup> cas, ceci nous donnant :

$$\begin{array}{r}
 \boxed{-} \quad 0 \quad 3 \quad 2 \\
 - \quad \quad 1 \quad 8 \quad 9 \\
 \hline
 \bullet \quad 9 \quad 6 \quad 8 \\
 + \quad \bullet \quad 8 \quad 1 \quad 1 \\
 \hline
 = \quad 1 \quad 7 \quad 7 \quad 9 \\
 \hline
 \boxed{-} \quad 2 \quad 2 \quad 1
 \end{array}$$

Nous avons donc envie de dire que  $n - g \in \{0; 1\}$  donne le nombre de compléments 9 plus un supplémentaires à faire ainsi que le nombre de signe moins à ajouter après avoir lu le résultat en ignorant le chiffre tout à gauche. Que c'est beau !

Les deux cas suivants achèvent notre exploration expérimentale sans contredire la conjecture.

$$\begin{array}{r}
 \quad \quad 0 \quad 3 \quad 2 \\
 + \quad \quad 1 \quad 8 \quad 9 \\
 \hline
 = \quad \quad 2 \quad 2 \quad 1 \\
 \hline
 \boxed{+} \quad 2 \quad 2 \quad 1
 \end{array}$$

$$\begin{array}{r}
 \quad \quad \boxed{-} \quad 0 \quad 3 \quad 2 \\
 + \quad \quad 1 \quad 8 \quad 9 \\
 \hline
 \bullet \quad 9 \quad 6 \quad 8 \\
 + \quad \quad 1 \quad 8 \quad 9 \\
 \hline
 = \quad 1 \quad 1 \quad 5 \quad 7 \\
 \hline
 \boxed{+} \quad 1 \quad 5 \quad 7
 \end{array}$$

**Exercice 1.** Dans la méthode ci-dessous, il existe des cas problématiques. Lesquels ?

**Exercice 2.** En laissant les cas problématiques de côté, démontrer le caractère général de la méthode que nous avons juste exposée via quelques exemples.

**2.3. Pourquoi ça marche ?** Revenons à nos entiers signés écrits sur 8 bits en reprenant le calcul de  $84 - 101 = -17$  tout en nous inspirant de ce qui a été vu dans la section précédente avec la base 10. Ici  $256 = 2^8$  joue le même rôle que  $1000 = 10^3$  avant. Nous avons ajouté des cases grises pour le bit du calcul intermédiaire, et nous avons indiqué les complémentations<sup>2</sup> en mettant sur fond vert les bits les plus à gauche..

$$\begin{array}{r}
 \begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & \text{gris} & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ \hline 0 & \text{gris} & 1 & 1 & 0 & 0 & 1 & 0 & 1 \\ \hline \end{array} \\
 - \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & \text{gris} & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ \hline 1 & \text{gris} & 0 & 0 & 1 & 1 & 0 & 1 & 1 \\ \hline \end{array} \\
 + \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 & 1 \\ \hline \end{array} \\
 = \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & \text{gris} & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ \hline \end{array} \\
 -
 \end{array}$$

2. Nous utiliserons « *complémentation* » comme abréviation de « *complément à 2 plus un* ».

La dernière complémentation avec ajout d'un signe vient du 0 sur fond gris et du fait qu'une seule complémentation préparatoire a été faite avant. Tout s'éclaire !

Reprenons de même le cas de  $84 - 61 = 23$  où l'on avait une retenue à ignorer. On fait comme précédemment.

$$\begin{array}{r}
 \begin{array}{|c|c|c|c|c|c|c|c|c|}
 \hline
 0 & \text{gris} & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 0 & \text{gris} & 0 & 1 & 1 & 1 & 1 & 0 & 1 \\
 \hline
 \end{array} \\
 - \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|c|}
 \hline
 0 & \text{gris} & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 1 & \text{gris} & 1 & 0 & 0 & 0 & 0 & 1 & 1 \\
 \hline
 \end{array} \\
 + \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|c|}
 \hline
 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 \end{array} \\
 = \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|c|}
 \hline
 0 & \text{gris} & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 \end{array} \\
 + \\
 \hline
 \end{array}$$

Avec le nouvel éclairage, nous n'avons qu'à garder les bits sur fond blanc, le signe du résultat étant positif<sup>3</sup>.

Dans les deux cas précédents, au lieu de considérer un bit intermédiaire, il suffit d'effectuer directement le calcul de retenue sur le bit de gauche modulo 2, ce calcul traduisant celui sur le nombre de complémentations préparatoires et du chiffre intermédiaire en gris dans nos deux exemples ci-dessus. On peut ainsi faire plus directement comme suit.

$$\begin{array}{r}
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 \\
 \hline
 \end{array} \\
 - \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 1 & 0 & 0 & 1 & 1 & 0 & 1 & 1 \\
 \hline
 \end{array} \\
 + \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 1 & 1 & 1 & 0 & 1 & 1 & 1 & 1 \\
 \hline
 \end{array} \\
 = \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
 \hline
 \end{array} \\
 - \\
 \hline
 \end{array}
 \qquad
 \begin{array}{r}
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 0 & 0 & 1 & 1 & 1 & 1 & 0 & 1 \\
 \hline
 \end{array} \\
 - \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 \hline
 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 \\
 \hline
 \end{array} \\
 + \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 \end{array} \\
 = \\
 \hline
 \begin{array}{|c|c|c|c|c|c|c|c|}
 \hline
 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\
 \hline
 \end{array} \\
 + \\
 \hline
 \end{array}$$

Plus généralement, avec la façon de stocker les entiers signés, nous avons alors les correspondances suivantes où les entiers additionnés  $a$  et  $b$  sont dans  $[-128; 127]$  et leur somme aussi, ce qui revient à ne considérer que les cas de non dépassement de capacité.

$$\begin{array}{cccc}
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 0 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 0 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 0 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 0 & \dots \\
 \hline
 \end{array} \\
 - & - & - & - \\
 \hline
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 1 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 1 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 1 & \dots \\
 \hline
 1 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 1 & \dots \\
 \hline
 1 & \dots \\
 \hline
 \end{array} \\
 + & + & + & + \\
 \hline
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 1 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 0 & \dots \\
 \hline
 \end{array} &
 \begin{array}{|c|c|}
 \hline
 1 & \dots \\
 \hline
 \end{array} \\
 = & = & = & = \\
 \text{Soustraction 1} & \text{Soustraction 2} & \text{Addition 1 de} & \text{Addition 2 de} \\
 \text{de deux naturels} & \text{de deux naturels} & \text{deux relatifs négatifs} & \text{deux relatifs négatifs}
 \end{array}$$

Les soustractions et l'addition 2 ne posent aucun souci. Par contre l'addition 1 est problématique avec son résultat positif. En fait cette addition contredit notre hypothèse de non dépassement de capacité. Nous allons voir pourquoi.

3. Il y a eu une seule complémentation préparatoire et le bit intermédiaire est 1.

Le cas de l'addition 1 correspond à  $(a; b) \in [-128; -1]^2$  tel que  $(128 + a) + (128 + b) \in [0; 127]$  soit  $0 \leq 256 + a + b \leq 127$  i.e.  $-256 \leq a + b \leq -129$  ce qui correspond à un dépassement de capacité comme annoncé.

Ceci achève de démontrer la validité des procédés d'addition et de soustraction d'entiers signés dans les cas de non dépassement de capacité. Notez que le cas évident d'une addition de deux naturels a été omis, et aussi que  $-a + b = b - a = b + (-a)$  et le fait qu'un changement de signe n'est autre qu'un complément à 1 plus 1 permettent de compléter les cas non indiqués ci-dessus.

**Exercice 3.** *Étudiez plus généralement le cas d'une base  $b \in \mathbb{N}_{\geq 3}$  quelconque.*

### 3. DIVISION PAR 2 ET DÉCALAGE À DROITE

**3.1. Entiers non signés.** Lorsque l'on calcule  $(n \div 2)$  le quotient entier par défaut d'un naturel non signé  $n$  divisé par 2, il suffit de retirer son bit de poids faible. Voici un exemple.

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 \\ \hline \end{array} \leftarrow n_0 = 128 + 64 + 16 + 4 + 2 + 1 = 215$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 1 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \leftarrow n_1 = n_0 \div 2 = 107$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 1 & 0 & 1 & 0 & 1 \\ \hline \end{array} \leftarrow n_2 = n_1 \div 2 = 53$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ \hline \end{array} \leftarrow n_3 = n_2 \div 2 = 26$$

Chaque étape correspond à un décalage vers la droite tout en complétant par un zéro à gauche.

**3.2. Entiers signés.** Examinons ce qu'il se passe lorsque l'on calcule  $(n \div 2)$  le quotient entier par défaut d'un naturel signé  $n < 0$  divisé par 2.

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 \\ \hline \end{array} \leftarrow n_0 = -128 + 64 + 16 + 4 + 2 + 1 = -41$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \leftarrow n_1 = n_0 \div 2 = -21 = -128 + 64 + 32 + 8 + 2 + 1$$

On en déduit que quelque soit le signe de l'entier signé, le quotient entier par défaut correspond à un décalage vers la droite tout en conservant le bit de signe tout à gauche. Notez bien que l'on décale tous les bits y compris celui du signe. Démontrons nos affirmations.

Le cas  $n \in [0; 127]$  étant évident, considérons  $n \in [-128; -1]$ . La représentation de  $n$  s'obtient alors via  $n = -2^7 + \sum_{0 \leq k \leq 6} b_k 2^k$  avec chaque  $b_k$  dans  $\{0; 1\}$ . Nous avons alors :

$$n \div 2 = -2^6 + \sum_{1 \leq k \leq 6} b_k 2^{k-1}$$

$$n \div 2 = -2 \times 2^6 + 2^6 + \sum_{0 \leq k \leq 5} b_{k-1} 2^k$$

$$n \div 2 = -2^7 + 2^6 + \sum_{0 \leq k \leq 5} b_{k-1} 2^k$$

La dernière identité correspond bien au décalage avec conservation de signe qui a été annoncé plus haut.

**3.3. Synthèse.** Nous constatons que nous avons deux types de décalage à droite. Dans le cas des entiers non signés, le décalage vers la droite est noté  $\ggg$  et il est appelé « *décalage logique à droite* ». Dans le cas des entiers signés, le décalage vers la droite est noté  $\gg$  et il est appelé « *décalage arithmétique à droite* ».

## 4. MULTIPLICATION PAR 2 ET DÉCALAGE À GAUCHE

4.1. **Entiers non signés.** Lorsque l'on calcule  $2n$  le double d'un naturel non signé  $n$ , il suffit de retirer son bit de poids fort tout en ajoutant un zéro comme nouveau bit de poids faible. Voici un exemple sans dépassement de capacité.

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ \hline \end{array} \leftarrow n_0 = 16 + 4 + 2 + 1 = 23$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 \\ \hline \end{array} \leftarrow n_1 = 2n_0 = 46$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ \hline \end{array} \leftarrow n_2 = 2n_1 = 92$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ \hline \end{array} \leftarrow n_3 = 2n_2 = 184$$

Chaque étape correspond à un décalage vers la gauche tout en complétant par un zéro à droite.

4.2. **Entiers signés.** Examinons ce qu'il se passe lorsque l'on calcule  $2n$  si  $n < 0$  est un naturel signé. Voici un exemple sans dépassement de capacité.

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 \\ \hline \end{array} \leftarrow n_0 = -128 + 64 + 16 + 4 + 2 + 1 = -41$$

$$\begin{array}{|c|c|c|c|c|c|c|c|} \hline 1 & 0 & 1 & 0 & 1 & 1 & 1 & 0 \\ \hline \end{array} \leftarrow n_1 = 2n_0 = -82 = -128 + 32 + 8 + 4 + 2$$

On effectue un décalage à gauche de tous les bits, y compris celui du signe, comme avec les entiers non signés. Mais est-ce toujours vrai? Que se passe-t-il si le 2<sup>e</sup> bit à gauche est nul? En fait, sous l'hypothèse de non dépassement de capacité, tout fonctionne sans souci comme nous allons l'expliquer proprement tout de suite.

Le cas  $n \in [0; 127]$  étant évident, considérons  $n \in [-128; -1]$ . La représentation de  $n$  s'obtient alors via  $n = -2^7 + \sum_{0 \leq k \leq 6} b_k 2^k$  avec chaque  $b_k$  dans  $\{0; 1\}$ . Nous avons alors :

$$2n = -2^8 + \sum_{0 \leq k \leq 6} b_k 2^{k+1}$$

$$2n = -2^8 + b_6 2^7 + \sum_{1 \leq k \leq 6} b_{k-1} 2^k$$

$$2n = (b_6 - 2) 2^7 + \sum_{1 \leq k \leq 6} b_{k-1} 2^k$$

Si  $b_6 = 1$ , nous avons bien le décalage annoncé. Supposons donc que  $b_6 = 0$ . Dans ce cas, en notant  $n = -2^7 + S$ , nous obtenons :

$$S = \sum_{0 \leq k \leq 6} b_k 2^k$$

$$S = \sum_{0 \leq k \leq 5} b_k 2^k$$

$$S \leq \sum_{0 \leq k \leq 5} 2^k$$

$$S \leq 2^6 - 1$$

Nous en déduisons que  $n = -2^7 + S \leq -2^7 + 2^6 - 1 = -2^6 - 1$ , puis ensuite  $2n \leq -2^7 - 2 < -2^7$  ce qui est bien un cas de dépassement de capacité.

**4.3. Synthèse.** Dans le cas des entiers signés ou non, nous utilisons le même décalage vers la gauche qui est noté  $\ll$ . Pour préciser sur quel type d'entiers on effectue un décalage, on parlera de « *décalage logique à gauche* » pour les entiers non signés, et de « *décalage arithmétique à gauche* » pour les entiers signés.