# MI-ARI
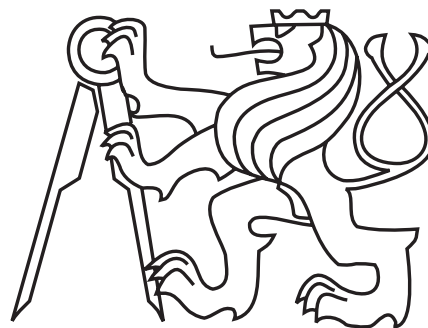## (Computer arithmetics)
## winter semester 2017/18

# CS. Number systems and basic operations

© **Alois Pluháček Pavel Kubalík**, 2017
**Department of digital design**
**Faculty of Information technology**
**Czech Technical University in Prague**

# CS. Number systems and basic operations

- **Number systems**
- **Number formats**
  - **Basic operation in general format**
- **Binary adders**
- **Format respecting addition**
- **Subtractor**
- **Subtractor using adder**
- **Signed number representation**
  - **Sign and magnitude**
  - **Radix complement**
  - **Diminished radix complement**
  - **Biased code**
    * **Type 0 of biased code**
    * **Type 1 of biased code**

# Number systems

number systems
- **position number systems**
  - **standard**
    - ∗ **decimal**
    - ∗ **binary (or dyadic)**
    - ∗ **octal**
    - ∗ **hexadecimal**
    - ∗ **ternary (or triadic)**
      etc.
  - **non-standard number systems**
    - ∗ **negative radix number systems** (Polish system)
    - ∗ **signed digit number systems**
    - ∗ **multiple radix number systems**
      **and others.**
- **non-position number systems**
  - **so caller Roman numerals**
  - **residue number systems (RNS)**, (Czech system)
    etc.

# Standard number system –position number system

$$A \sim a_n a_{n-1} \ldots a_0, a_{-1} \ldots a_{-m}$$

$$A = \sum_{i=-m}^{n} a_i z^i$$

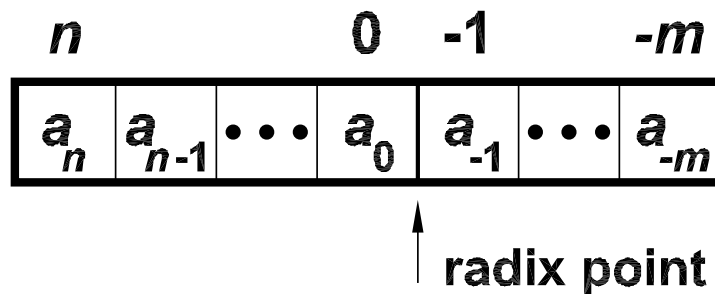$$z \geq 2$$

radix (or base) of system — natural number

$$a_i \in \langle 0; z)$$

digits — non-negative integers

$$0 \leq A \leq z^{n+1} - z^{-m}$$

!!! Negative numbers can not be represented !!!

# Number formats



| | | | | | | |
|---|---|---|---|---|---|---|
| $n$ | | | $0$ | $-1$ | | $-m$ |

$a_n$ $a_{n-1}$ $\cdots$ $a_0$ $a_{-1}$ $\cdots$ $a_{-m}$

↑ radix point

$n$ ... **highest position**

$-m$ ... **lowest position**

$A \sim \mathbf{a_n a_{n-1} \ldots a_0, a_{-1} \ldots a_{-m}}$

$A = a_n z^n + a_{n-1} z^{n-1} + \ldots + a_0 + a_{-1} z^{-1} \ldots a_{-m} z^{-m}$

$z$ ... **radix** (or base) of the number system

$\mathcal{Z} = z^{n+1}$ **module of the format** — it is out of format

$\varepsilon = z^{-m}$ **unit of the format** — smallest positive number

in the format

**representable numbers** $A$: $\boxed{0 \leq A = k \cdot \varepsilon < \mathcal{Z}}$ ,

$k$ **is an integer**

$k = A/\varepsilon = A^* \implies \boxed{A = A^* \text{ of units } \varepsilon}$

$A^*$ ☞ **digit number at notation of** $A$

$\varepsilon$ ☞ **radix point position**

# Basic operation in general format

## Addition and subtraction

### the same format of both operands and the result:

$$\left.\begin{array}{l} A = A^* \cdot \varepsilon \\ B = B^* \cdot \varepsilon \end{array}\right\} \;\Rightarrow\; A \pm B = A^* \cdot \varepsilon \pm B^* \cdot \varepsilon = (A^* \pm B^*) \cdot \varepsilon$$

$$\text{Ex.:}\quad z = 10, \;\; \mathcal{Z} = 10 = 10^{n+1}, \;\; \varepsilon = 0,01 = 10^{-m}$$

$$n = 0, \;\; m = 2 \;\; (\text{or } -m = -2)$$

$$A = 1,23 \;\;\Rightarrow\;\; A^* = 1,23/0,01 = 123$$

$$B = 4,56 \;\;\Rightarrow\;\; B^* = 4,56/0,01 = 456$$

$$A = 1,23 + 4,56 = (123 + 456) \cdot 0,01 =$$
$$= 579 \; \cdot 0,01 = 5,79$$

### different formats:

**transformation numbers into suitable format —** zeroes adding

$$\text{Ex.:}\quad \textbf{1,234+56,7} \;=\; \textbf{01,234+56,700} \;=\; \textbf{57,934}$$

---

**Conclusion:** Addition a subtraction (in general format) **can be easy transformed on addition a subtraction of integers.**

---

## Multiplication:

$$\left. \begin{array}{l} A = A^* \cdot \varepsilon_A \\ B = B^* \cdot \varepsilon_B \end{array} \right\} \quad \Rightarrow \quad \begin{array}{l} A \cdot B = A^* \cdot \varepsilon_A \cdot B^* \cdot \varepsilon_B = \\ = (A^* \cdot B^*) \cdot \varepsilon_A \cdot \varepsilon_B \end{array}$$

**Ex.:** $z = 10$

$\mathcal{Z}_A = 10, \quad \varepsilon_A = 0,01, \ n_A = 0, \ m_A = 2$

$\mathcal{Z}_B = 100, \ \varepsilon_B = 0,1, \quad n_B = 1, \ m_B = 1$

$7,01 \cdot 80,3 = (701 \cdot 803) \cdot 0,001 =$

$\qquad\qquad = 562\,903 \cdot 0,001 = 562,903$

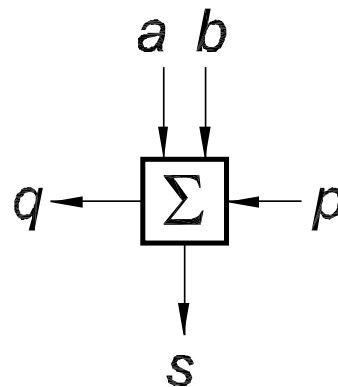**Conclusion: Multiplication**(in general format) **can be easy transformed on addition a subtraction of integers.**

**Full-adder**

| $a$ | $b$ | $p$ | $q$ | $s$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

$$s = a \oplus b \oplus p =$$
$$= \overline{a}\,\overline{b}p + \overline{a}b\overline{p} + a\overline{b}\,\overline{p} + abp$$

$$q = \mathsf{M}_3(a, b, p) =$$
$$= ab + ap + bp =$$
$$= ab \oplus ap \oplus bp =$$
$$= ab + (ap \oplus bp)$$

**Half-adder**

| $a$ | $b$ | $q$ | $s$ |
|-----|-----|-----|-----|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |

$$s = a \oplus b$$
$$= \overline{a}b + a\overline{b}$$
$$q = a \cdot b$$





**CS – 7** © A. Pluháček / P. Kubalík 2017

**Ripple-carry adder**
(also *parallel adder*)

$n = 3$
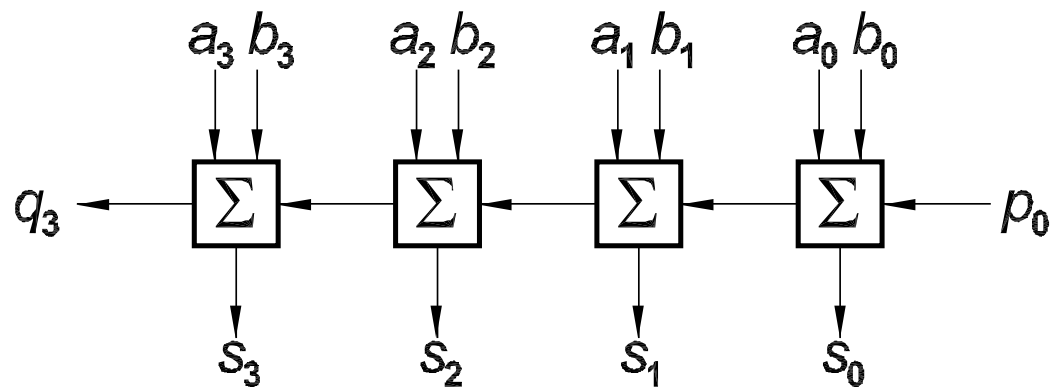$\mathcal{Z} = 16$
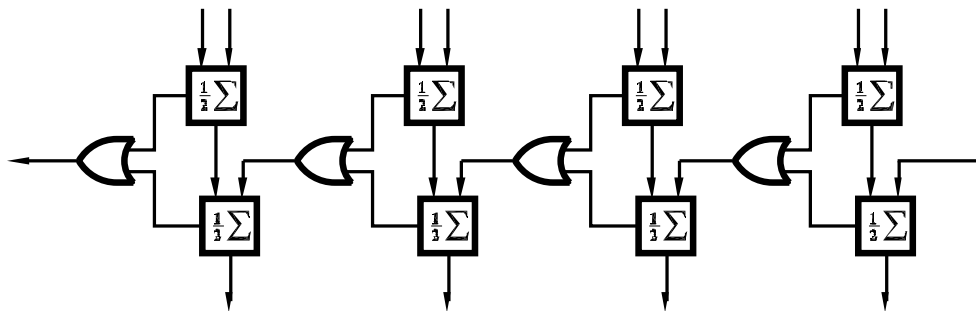$A \sim a_3 a_2 a_1 a_0$
$B \sim b_3 b_2 b_1 b_0$
$S \sim s_3 s_2 s_1 s_0$
$p_{i+1} = q_i$



$$S = A + B + p_0 - q_n \cdot \mathcal{Z}$$

the same using half-adders

# Format respecting addition

**Output of adder:** $S = A + B + p_0 - q_n \cdot \mathcal{Z}$

**Let** $p_0 = 0$
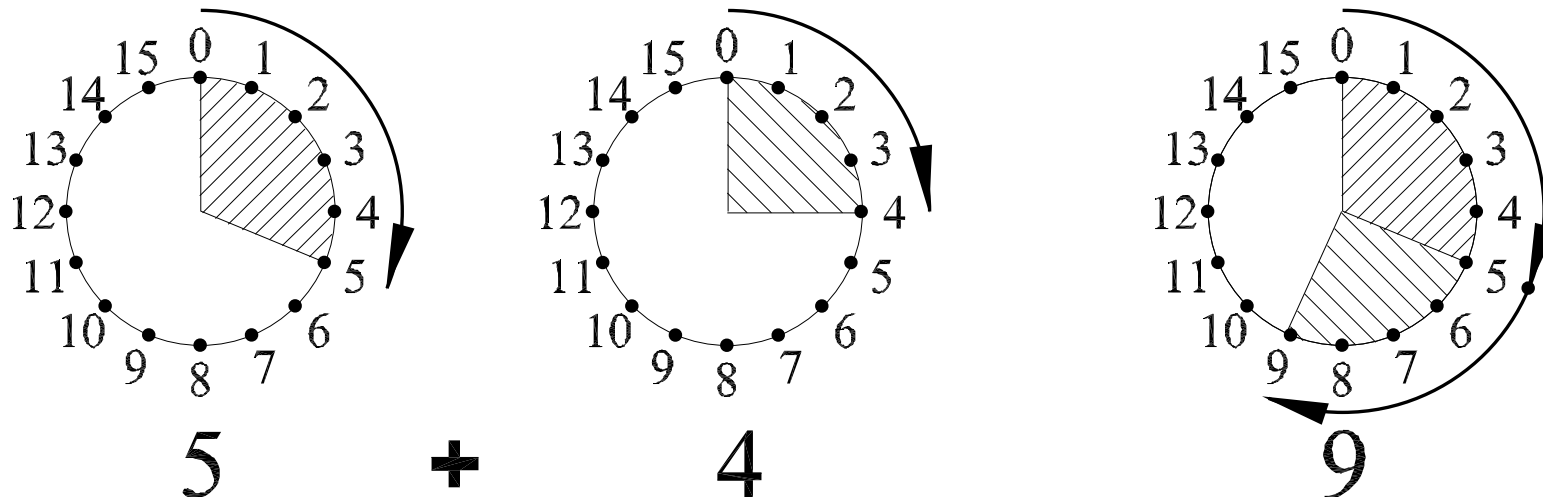
**(or half adder is used in zero order of adder):**

$$\boxed{S = A + B - q_n \cdot \mathcal{Z}}$$

$S$ **differ from** $A + B$ **by multiple of** $\mathcal{Z}$ **so that**    $\mapsto$
$S \equiv A + B \pmod{\mathcal{Z}}$

**graphic view (analogy of clock face):**

$$0101 + 0100 = {}_0\mathbf{1001} \rightarrow 1001$$



$$5 \qquad + \qquad 4 \qquad\qquad\qquad 9$$

$$0101 + 1110 = {}_1 \mathbf{0011} \rightarrow \mathbf{0011}$$



$$5 \qquad + \qquad 14 \qquad\qquad 3$$

**passing thru** $\quad \Leftrightarrow \quad$ **carry from the highest order** $q_n = 1$

**in this case (additon using unsigned numbers):**

$$q_n = 1 \quad \Rightarrow \quad A + B \geq \mathcal{Z} \qquad (\mathcal{Z} = 1\,0000_2 = 16_{10})$$

$$q_n = 1 \quad \Rightarrow \quad \text{overflow} \sim \textbf{the result is out of format}$$

**however generally:** **!!! carry $\not\equiv$ overflow !!!**

# Subtractor

**Full adder (!** *How to do it in wrong way.* **!)**

| $a$ | $b$ | $v$ | $u$ | $r$ |
|-----|-----|-----|-----|-----|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 |

$$r = a \oplus b \oplus v =$$
$$= \overline{a}\,\overline{b}v + \overline{a}b\overline{v} + a\overline{b}\,\overline{v} + abv$$
$$u = = \overline{a}b + \overline{a}v + bv$$

$v_i$     **borrow for order** $i$

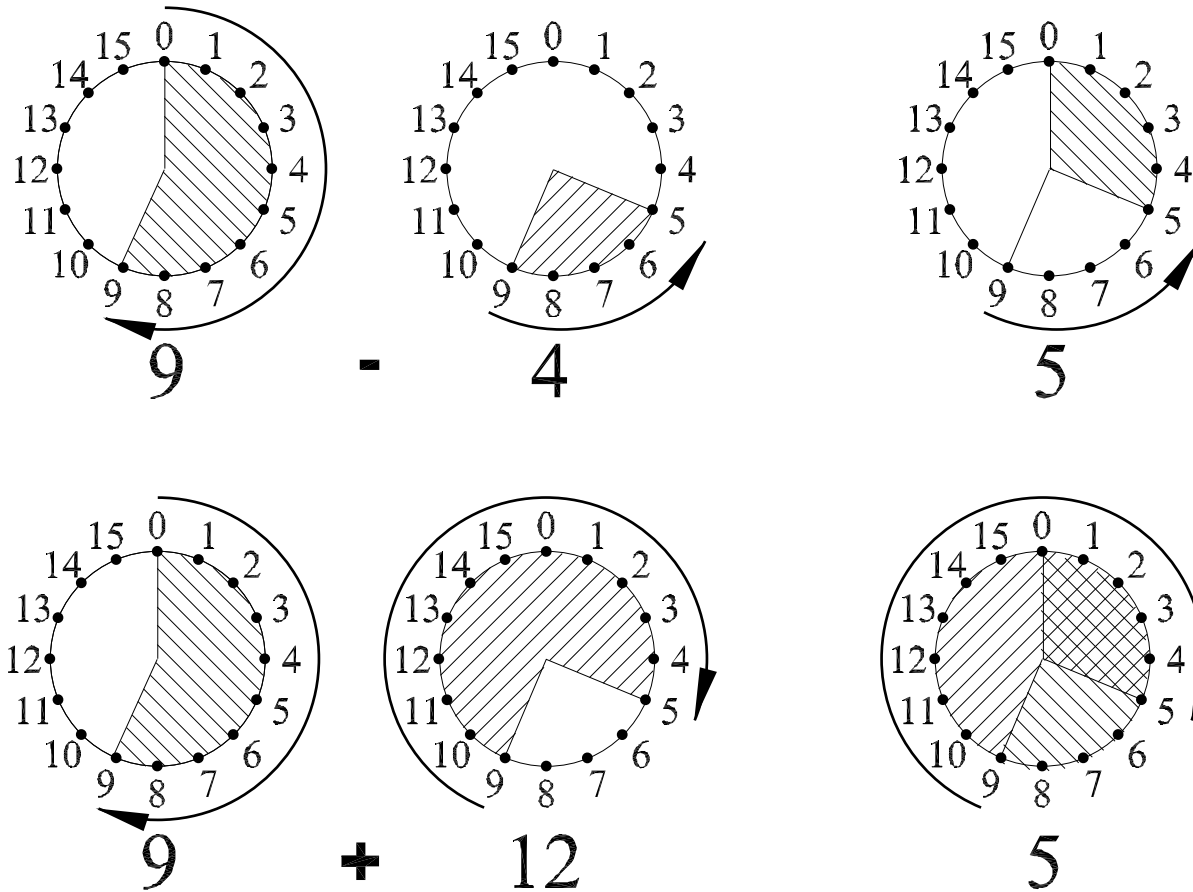$u_i$     **borrof from order** $i$

**etc. - similarly as for addition**

**Is it possible modify an adder for subtraction?**

$$A - B \equiv A + (\mathcal{Z} - B) \quad (\text{mod } \mathcal{Z})$$

**How to find $\mathcal{Z} - B$?**

$$X = \sum_{i=0}^{n} x_i z^i \qquad x_i \in \langle 0, z-1 \rangle$$

$$X_{max} = \sum_{i=0}^{n} (z-1)z^i = \sum_{j=1}^{n+1} z^j - \sum_{i=0}^{n} z^i = z^{n+1} - 1 \; = $$

$$= \; \mathcal{Z} - 1$$

$z = 2$:  $\quad X_{max} = \boxed{11\ldots 11} = \mathcal{Z} - 1$

$\boxed{\text{xx}\ldots\text{xx}}$ — module format

$\mathcal{Z} = \boxed{11\ldots 11} + 1$

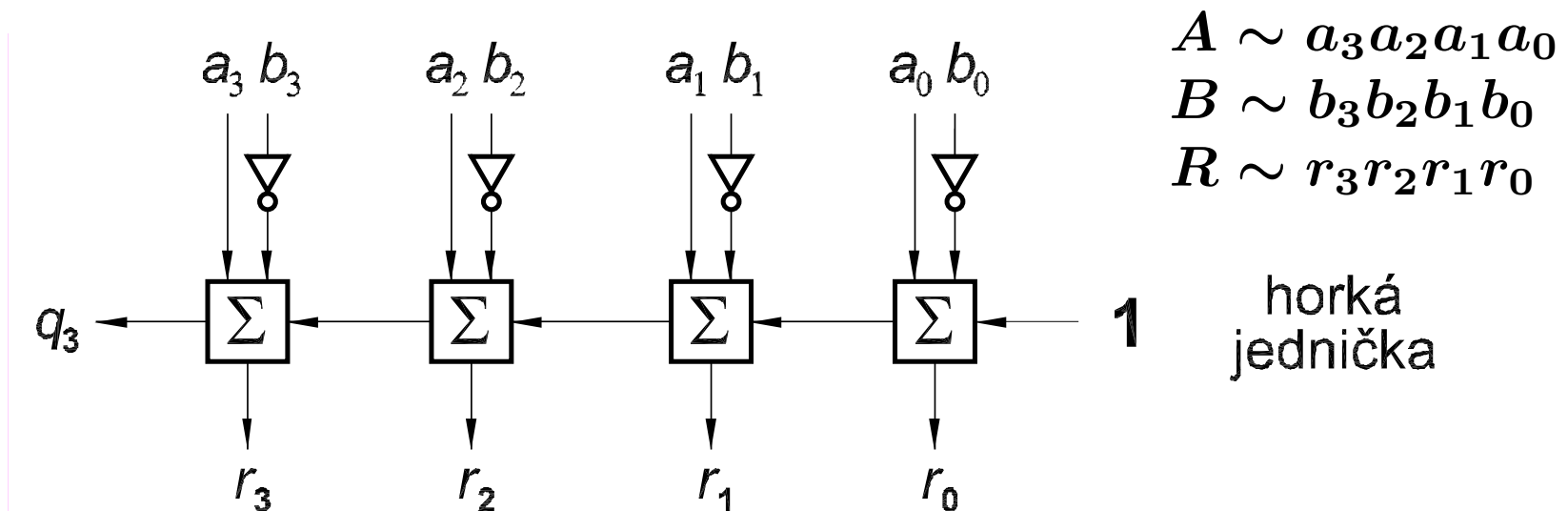$\mathcal{Z} - B = \boxed{11\ldots 11} - B + 1$

$\boxed{\boxed{\mathcal{Z} - B = \overline{B} + 1}}$  $\qquad \overline{B} \ldots$ **negation of all bits**

$\ldots + 1$ — **so called** hot one

**note.:**  $B = 0 \Rightarrow \mathcal{Z} - B = \mathcal{Z} \equiv 0 \pmod{\mathcal{Z}}$

$\overline{B} + 1 = \boxed{11\ldots 11} + 1 = {}_1\boxed{00\ldots 00}$

$$A \sim a_3 a_2 a_1 a_0$$
$$B \sim b_3 b_2 b_1 b_0$$
$$R \sim r_3 r_2 r_1 r_0$$



horká jednička

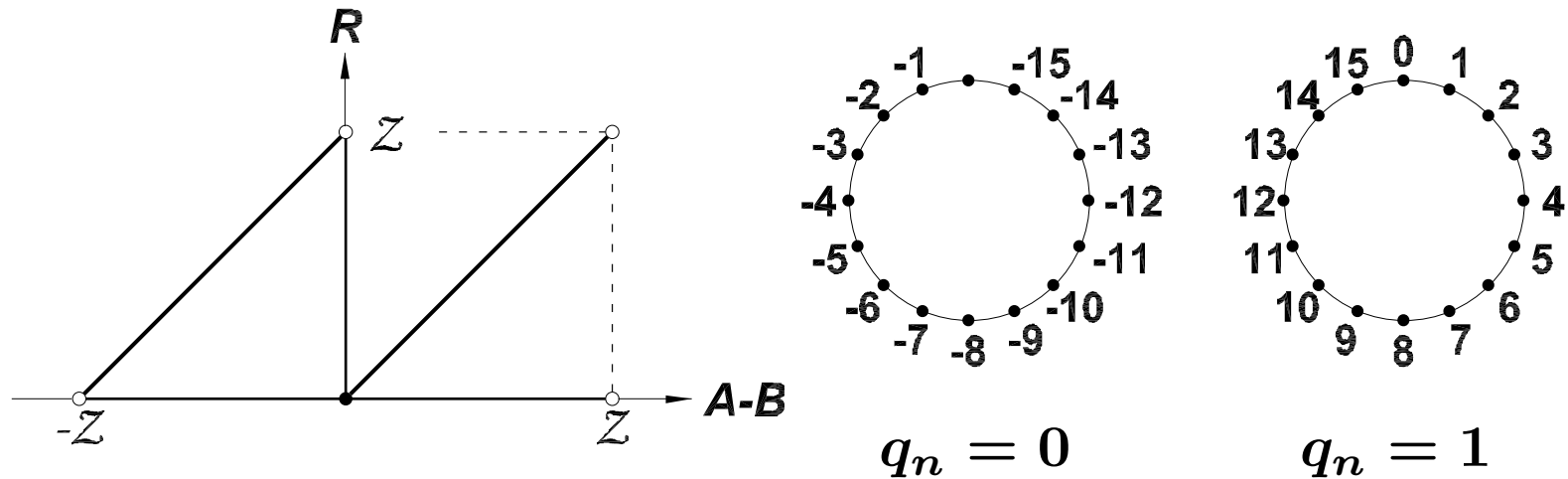$$R = A + (\mathcal{Z} - B) - q_n \cdot \mathcal{Z} = A - B + (1 - q_n) \cdot \mathcal{Z}$$

$$R = A - B + \overline{q_n} \cdot \mathcal{Z} \qquad 0 \leq R < \mathcal{Z}$$

$$q_n = 1 \quad \Rightarrow \quad R = A - B \geq 0$$
$$q_n = 0 \quad \Rightarrow \quad R = A - B + \mathcal{Z} < \mathcal{Z} \quad \Rightarrow A - B < 0$$

$$q_n = 1 \quad \Rightarrow \quad A \geq B \qquad R = A - B$$
$$q_n = 0 \quad \Rightarrow \quad A < B \qquad R = \mathcal{Z} - (B - A)$$

## complementary pseudocode



$q_n = 0$      $q_n = 1$

**If it is $q_n = 0$,  it means $B > A$,  then**

$$R = \mathcal{Z} - (B - A) \quad \Rightarrow \quad B - A = \mathcal{Z} - R$$

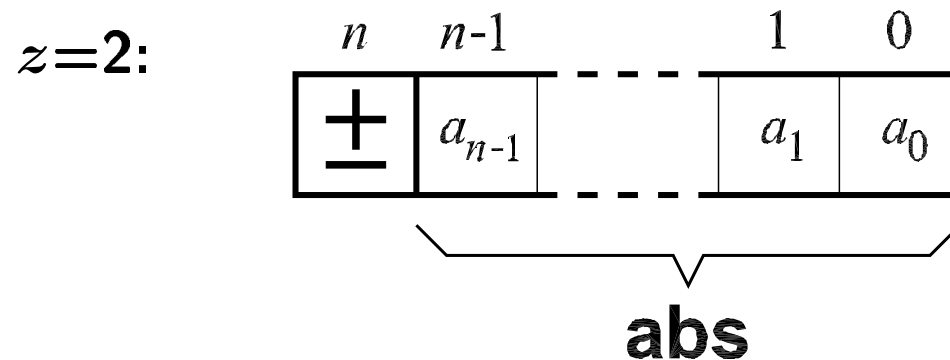$$\Rightarrow \quad \boxed{B - A = \overline{R} + 1}$$

# Signed number representation

## 5 ways to represent
## negative (as well as non-negative) number $X$:

1. **sign and magnitude** ... $\mathcal{P}(X)$

2. **radix complement** ... $\mathcal{D}(X)$
   **2's complement** - $z = 2$
   **10's complement** - $z = 10$

3. **diminished radix complement** ... $\mathcal{I}(X)$
   **1's complement** - $z = 2$
   **9's complement** - $z = 10$

4. **biased code** (or **excess $K$**) ... $\mathcal{A}(X)$
   a. **type 0** ... $\mathcal{A}_0(X)$
   b. **type 1** ... $\mathcal{A}_1(X)$

# Sign and magnitude
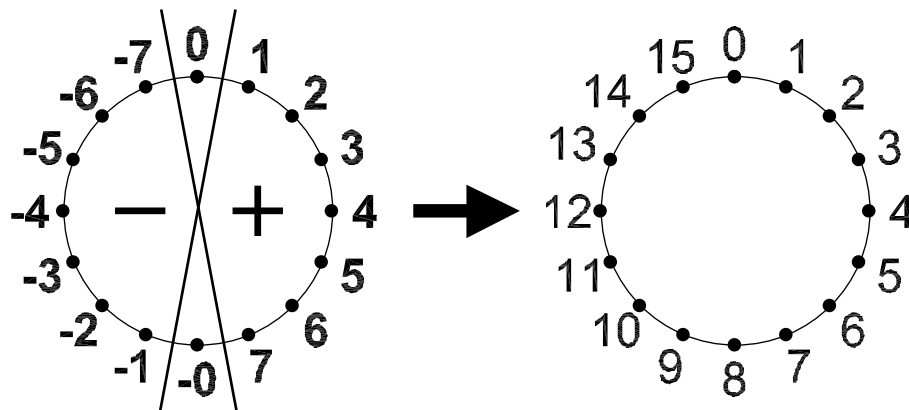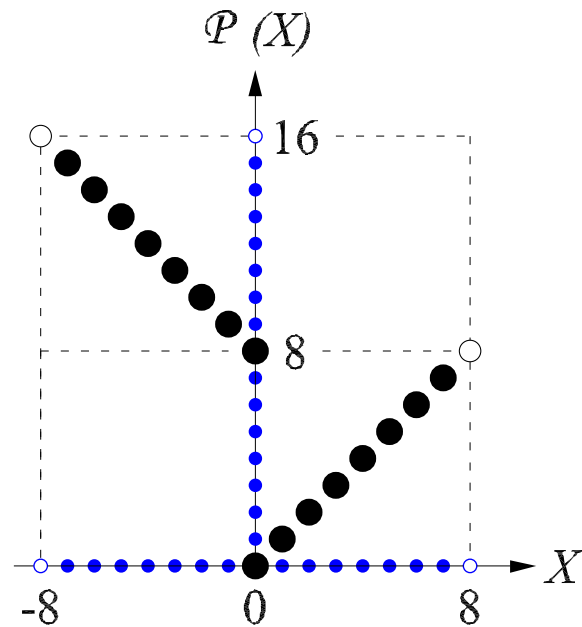
**sign & abs value (magnitude):**

$z=2$:



$$\text{sign bit} = \begin{cases} 0 & \text{for } X \geq 0 \\ 1 & \text{for } X \leq 0 \end{cases}$$

$$\mathcal{P}(X) = \begin{cases} X & \text{for } X \geq 0 \\ 2^n + |X| & \text{for } X \leq 0 \end{cases}$$

$-\frac{1}{2}\mathcal{Z} < X < \frac{1}{2}\mathcal{Z}$  ☞  **symmetric range**

**2 zero representation:** $\begin{cases} \text{so called „positive zero"} \quad 0\ 0\ldots 0 \\ \text{so called „negative zero"} \quad 1\ 0\ldots 0 \end{cases}$

| $X$ | $\mathcal{P}(X)$ |
|:---:|:---:|
| **0** | **0000** |
| **1** | **0001** |
| **2** | **0010** |
| **3** | **0011** |
| **4** | **0100** |
| **5** | **0101** |
| **6** | **0110** |
| **7** | **0111** |
| **-0** | **1000** |
| **-1** | **1001** |
| **-2** | **1010** |
| **-3** | **1011** |
| **-4** | **1100** |
| **-5** | **1101** |
| **-6** | **1110** |
| **-7** | **1111** |

**addition :**

$$ \boxed{(zA, aA) + (zB, aB) \rightarrow (zC, aC)} \ , $$

**where $zA$, $zB$ and $zC$ are the sign bits and $aA$, $aB$ and $aC$ are the magnitudes**

```
if (zA=zB)      {aA+aB → aC;
                 zA → zC;
                 if (q = 1) overflow; }
else            {aA+aB+1 → aC;
                 zA → zC;
                 if (q = 0)      {aC+1 → aC;
                                  zC → zC; } }
```

**$q$ carry-out from higher order**

**sign change :**

$$\mathcal{P}(-X) \to \mathcal{P}(X) \quad \Longleftarrow \quad \overline{\mathbf{MSB}} \to \mathbf{MSB}$$

**MSB - bit in higher order (first from left)**
**- Most Significant Bit**

**subtraction :**

$$A - B = A + (-B)$$

**to swap $\overline{zB}$ by $zB$**

**else alike the addition**

**absolute value :**

$$\mathcal{P}(X) \to \mathcal{P}(|X|) \quad \Longleftarrow \quad \mathbf{0} \to \mathbf{MSB}$$

# Radix complement

radix complement $\begin{cases} \textbf{2's complement for } z = 2 \\ \textbf{10's complement for } z = 10 \end{cases}$

$$\mathcal{D}(X) = \begin{cases} X & \textbf{pro } X \geq 0 \\ \mathcal{Z} + X = \mathcal{Z} - |X| & \textbf{pro } X < 0 \end{cases}$$
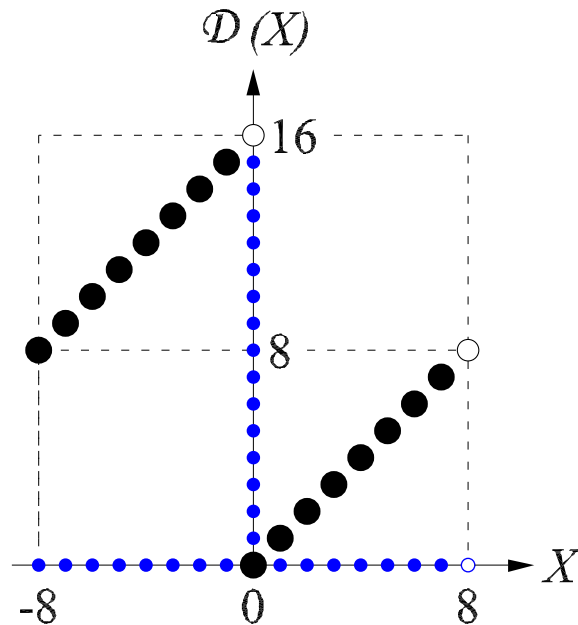
$$-\tfrac{1}{2}\mathcal{Z} \leq X < \tfrac{1}{2}\mathcal{Z} \qquad \text{☞} \qquad \textbf{asymmetric range}$$

$$\textbf{MSB} = 0 \quad \Longleftrightarrow \quad X \geq 0$$
$$\textbf{MSB} = 1 \quad \Longleftrightarrow \quad X < 0$$

$$\mathcal{D}(X) \equiv X \ (\textbf{mod } \mathcal{Z})$$

**z=2:**



| $X$ | $\mathcal{D}(X)$ |
|:---:|:---:|
| **0** | **0000** |
| **1** | **0001** |
| **2** | **0010** |
| **3** | **0011** |
| **4** | **0100** |
| **5** | **0101** |
| **6** | **0110** |
| **7** | **0111** |
| **-8** | **1000** |
| **-7** | **1001** |
| **-6** | **1010** |
| **-5** | **1011** |
| **-4** | **1100** |
| **-3** | **1101** |
| **-2** | **1110** |
| **-1** | **1111** |

**addition:** $\mathcal{D}(A+B) \equiv A+B \equiv \mathcal{D}(A) + \mathcal{D}(B) \quad (\textbf{mod } \mathcal{Z})$

☞
> **add up $\mathcal{D}(A) + \mathcal{D}(B)$ and**
> **ignore carry $q_n$ from higher order**

**subtraction:** $\mathcal{D}(A-B) \equiv A - B \equiv \mathcal{D}(A) - \mathcal{D}(B) \;(\textbf{mod } \mathcal{Z})$

☞
> **subtract $\mathcal{D}(B)$ from $\mathcal{D}(A)$ and**
> **ignore borrow $v_n$ from higher order**

**or**

☞
> **convert subtraction on addition, that is**
> **add up $\mathcal{D}(A) + \overline{\mathcal{D}(B)} + 1$ and**
> **ignore carry $q_n$ from higher order**

**Carry** (eventually borrow) **is ignored.**

**?**     **How to detect overflow?**
**(that is $A+B \geq \frac{1}{2}\mathcal{Z}$ or $A+B < -\frac{1}{2}\mathcal{Z}$)**     **?**

**overflow during addition in 2's complement code:**

**1.** $A < 0$ a $B \geq 0 \;\Rightarrow\; A \leq A + B < B$
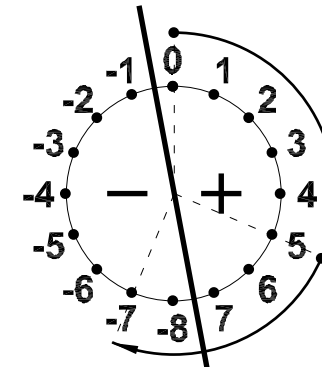
$B < 0$ a $A \geq 0 \;\Rightarrow\; B \leq A + B < A$

**!** <u>**in this case overflow can not occur**</u> **!**

**2.** $A \geq 0$ a $B \geq 0$

**overflow:** $A + B \geq \frac{1}{2}\mathcal{Z}$

**result has opposite sign**

**see example** $5 + 4 \rightarrow -7$

**3.** $A < 0$ a $B < 0$

**overflow:** $A + B < -\frac{1}{2}\mathcal{Z}$

**result has oposite sign**

**see example** $(-5) + (-5) \rightarrow 6$

**overflow during addition:**    **same sign of both operands and opposite sign of result**

| | | | |
|:-:|:-:|:-:|:-:|
| $+$ | $+$ | $\rightarrow$ | $-$ |
| $-$ | $-$ | $\rightarrow$ | $+$ |

**overflow during addition:**

                   **subtraction converted on addition**

$\Rightarrow$ **There is nothing to solve.**

**will be referred to:** $\mathcal{D}(A) + \mathcal{D}(B) - q_n \cdot \mathcal{Z} = S$

$$\mathcal{D}(A) \sim a_n^{\mathcal{D}} a_{n-1}^{\mathcal{D}} \ldots a_1^{\mathcal{D}} a_0^{\mathcal{D}}$$

$$\mathcal{D}(B) \sim b_n^{\mathcal{D}} b_{n-1}^{\mathcal{D}} \ldots b_1^{\mathcal{D}} b_0^{\mathcal{D}}$$

$$S \sim s_n^{\mathcal{D}} s_{n-1}^{\mathcal{D}} \ldots s_1^{\mathcal{D}} s_0^{\mathcal{D}}$$

$over$ — **overflow**

**detection of overflow:**

**1**  $a_n^{\mathcal{D}} = b_n^{\mathcal{D}} = 0$    **a**    $s_n^{\mathcal{D}} = 1$        **or**

$a_n^{\mathcal{D}} = b_n^{\mathcal{D}} = 1$    **a**    $s_n^{\mathcal{D}} = 0$

$$over = \overline{a_n^{\mathcal{D}}} \cdot \overline{b_n^{\mathcal{D}}} \cdot s_n^{\mathcal{D}} + a_n^{\mathcal{D}} \cdot b_n^{\mathcal{D}} \cdot \overline{s_n^{\mathcal{D}}}$$
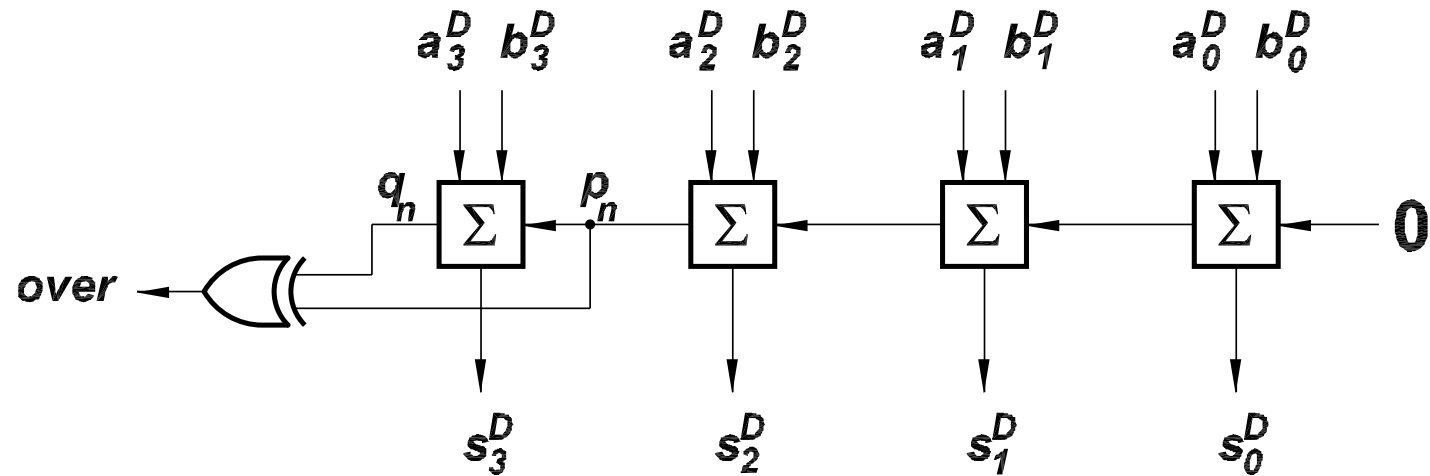
**2**

| $a_n^{\mathcal{D}}$ | $b_n^{\mathcal{D}}$ | $p_n^{\mathcal{D}}$ | $q_n^{\mathcal{D}}$ | $s_n^{\mathcal{D}}$ | |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | $q_n = p_n$ |
| ⊡0 | ⊡0 | 1 | 0 | ⊡1 | $q_n \neq p_n$ |
| 0 | 1 | 0 | 0 | 1 | $q_n = p_n$ |
| 0 | 1 | 1 | 1 | 0 | $q_n = p_n$ |
| 1 | 0 | 0 | 0 | 1 | $q_n = p_n$ |
| 1 | 0 | 1 | 1 | 0 | $q_n = p_n$ |
| ⊡1 | ⊡1 | 0 | 1 | ⊡0 | $q_n \neq p_n$ |
| 1 | 1 | 1 | 1 | 1 | $q_n = p_n$ |

$$over = q_n \oplus p_n$$

$z=2$:

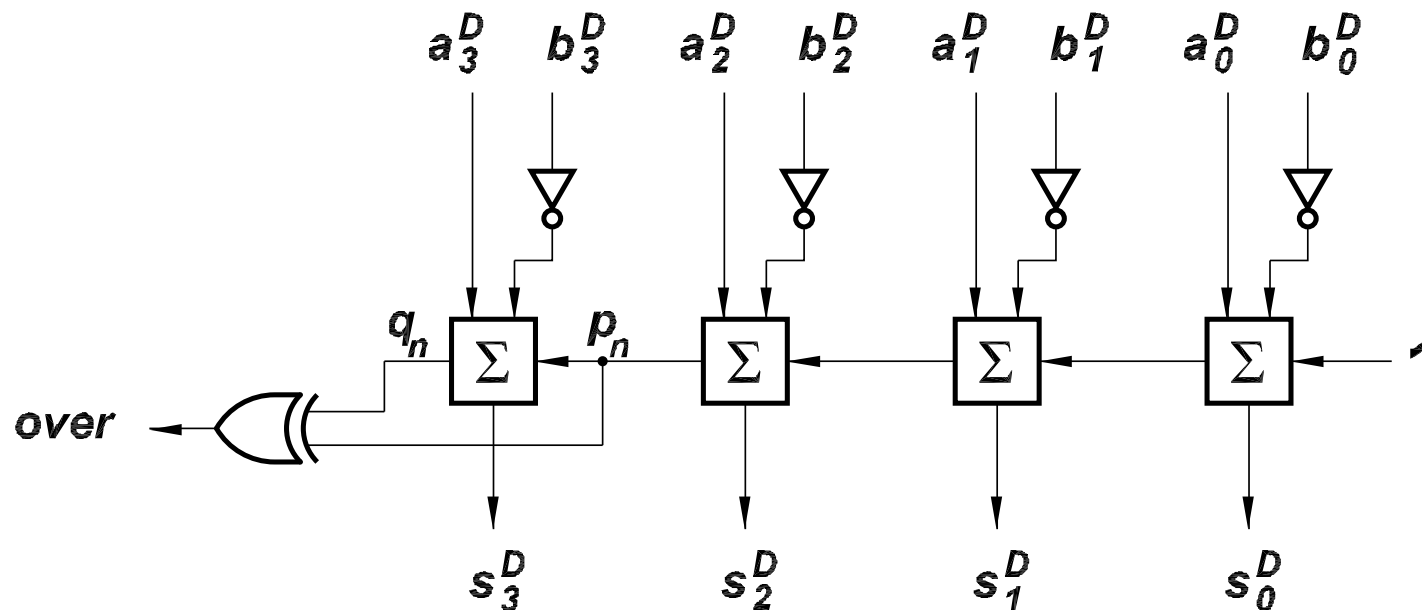**addition :** $\qquad \mathcal{D}(A + B) \equiv \mathcal{D}(A) + \mathcal{D}(B) \qquad (\textbf{mod } \mathcal{Z})$

$z=2$:

**sign change :**

$$\boxed{\mathcal{D}(-X) \equiv \overline{\mathcal{D}(X)}+1 \ (\textbf{mod } \mathcal{Z})}$$

**subtraction :** $A - B = A + (-B)$

# Diminished radix complement

radix complement $\begin{cases} \textbf{1's complement for } z = 2 \\ \textbf{9's complement for } z = 10 \end{cases}$

$$\mathcal{I}(X) = \begin{cases} X & \textbf{for } X \geq 0 \\ \overline{|X|} & \textbf{for } X \leq 0 \end{cases}$$

$-\frac{1}{2}\mathcal{Z} < X < \frac{1}{2}\mathcal{Z}$  ☞  **symmetric range**

**two zero representation:** $\begin{cases} \textbf{so called „positive zero" } 0\ 0\ldots0 \\ \textbf{so called „negative zero" } 1\ 1\ldots1 \end{cases}$

$$\begin{array}{l} \textbf{MSB} = 0 \implies X \geq 0 \\ \textbf{MSB} = 1 \implies X \leq 0 \end{array}$$

$T + \overline{T} = 11\ldots11 = \mathcal{Z}-1$      $\boxed{\mathcal{I}(X) \equiv X \pmod{\mathcal{Z}-1}}$

z=2:



| $X$ | $\mathcal{I}(X)$ |
|:---:|:---:|
| **0** | **0000** |
| **1** | **0001** |
| **2** | **0010** |
| **3** | **0011** |
| **4** | **0100** |
| **5** | **0101** |
| **6** | **0110** |
| **7** | **0111** |
| **-7** | **1000** |
| **-6** | **1001** |
| **-5** | **1010** |
| **-4** | **1011** |
| **-3** | **1100** |
| **-2** | **1101** |
| **-1** | **1110** |
| **- 0** | **1111** |

**addition :** $\qquad \mathcal{I}(A + B) \equiv \mathcal{I}(A) + \mathcal{I}(B) \qquad (\text{mod } \mathcal{Z}{-}1)$

$\mathcal{I}(A) + \mathcal{I}(B) \geq \mathcal{Z} \qquad \Rightarrow \qquad q{=}1$

**necessary to subtract $\mathcal{Z}$ and add 1 $\qquad \Longrightarrow$ circular carry**



**feedback $\qquad$ ☞ $\qquad$ sequential circuit**

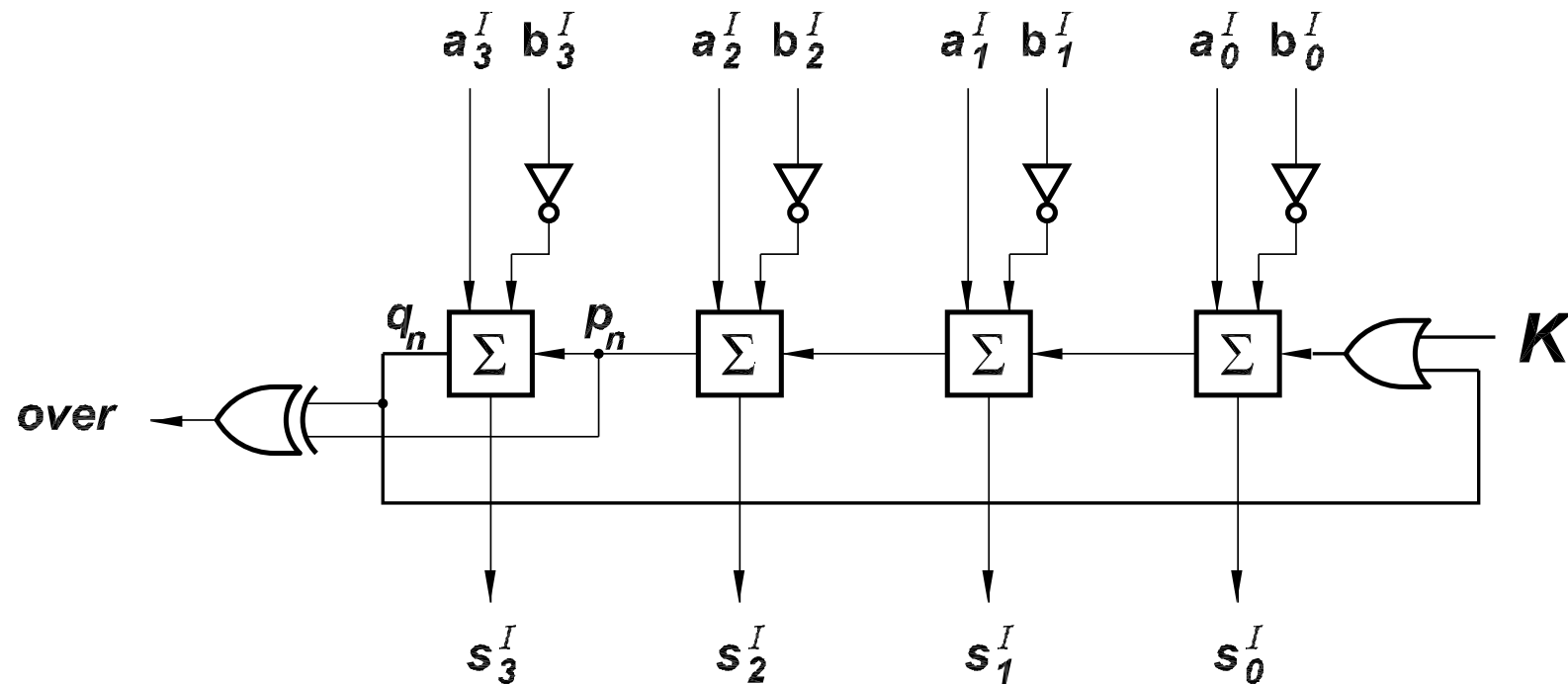**solution: $\quad$ correction $K$ on input carry into lower order**

$\qquad\qquad K = 1$, **in a case, when carry go over all orders**
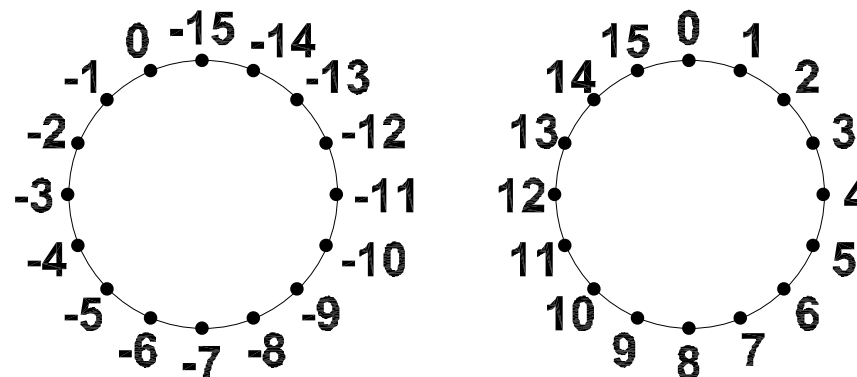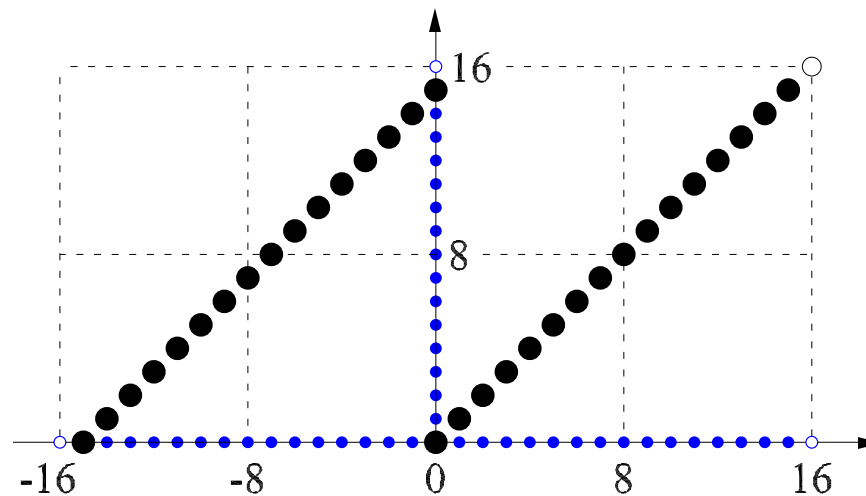
$z=2$:

**sign change :**

$$\boxed{\mathcal{I}(-X) \equiv \overline{\mathcal{I}(X)}}$$

**subtraction :** $A - B = A + (-B)$

**unsigned numbers + adder for radix complement code**    ⇨
⇨ **radix complement pseudo-code** (comparison
with 2's complement pseudo-code)

# Biased code

$$\boxed{\mathcal{A}(X) = X + K}\,,$$

where $K$ is „suitable" constant

„suitable" constants:

$$\frac{1}{2}\mathcal{Z} \qquad \text{☞} \qquad \mathcal{A}_0(X) \quad \ldots \quad \text{biased code type 0}$$

$$\frac{1}{2}\mathcal{Z}-1 \quad \text{☞} \quad \mathcal{A}_1(X) \quad \ldots \quad \text{biased code type 1}$$

**Code is monotonous — growing.**

**addition :** $\qquad \boxed{\mathcal{A}(A + B) = \mathcal{A}(A) + \mathcal{A}(B) - K}$

**subtraction :** $\boxed{\mathcal{A}(A - B) = \mathcal{A}(A) - \mathcal{A}(B) + K}$

**It is true (of course)** *iff overflow doesn't occurs !*

# Biased code - type 0

$$\boxed{\mathcal{A}_0(X) = X + \frac{1}{2}\mathcal{Z}}$$

$$-\frac{1}{2}\mathcal{Z} \leq X < \frac{1}{2}\mathcal{Z} \qquad \text{☞} \qquad \textbf{asymmetric range}$$

$$\boxed{\begin{aligned} \textbf{MSB} = 1 &\iff X \geq 0 \\ \textbf{MSB} = 0 &\iff X < 0 \end{aligned}}$$

$$\mathcal{A}_0(X) \equiv \mathcal{D}(X) \ (\textbf{mod } \tfrac{1}{2}\mathcal{Z})$$

$$\boxed{\textbf{! } \mathcal{A}_0(X) \textbf{ and } \mathcal{D}(X) \textbf{ differ in the MSB only !}}$$

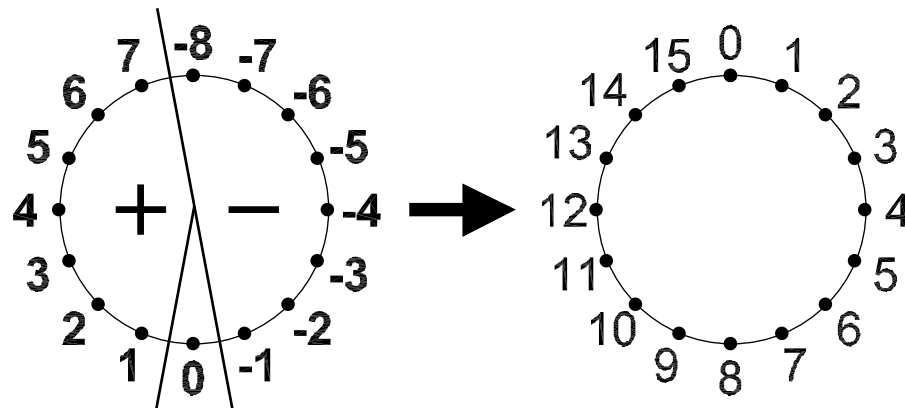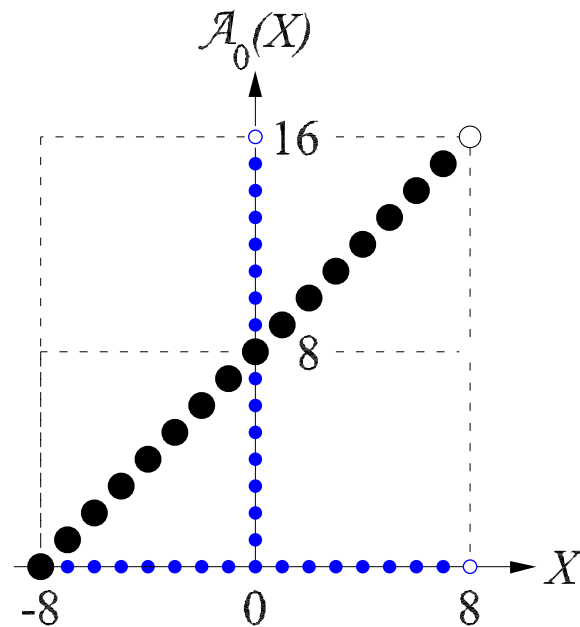$$\textbf{MSB}(\mathcal{A}_0(X)) = \overline{\textbf{MSB}(\mathcal{D}(X))}$$

$\implies$ **operations** (addition, subtraction, sign change etc.):

$$\mathcal{A}_0(operands) \to \mathcal{D}(operands)$$

$$\mathcal{D}(result) \to \mathcal{A}_0(result)$$

**z=2:**



| $X$ | $\mathcal{A}_0(X)$ |
|:---:|:---:|
| -8 | 0000 |
| -7 | 0001 |
| -6 | 0010 |
| -5 | 0011 |
| -4 | 0100 |
| -3 | 0101 |
| -2 | 0110 |
| -1 | 0111 |
| 0 | 1000 |
| 1 | 1001 |
| 2 | 1010 |
| 3 | 1011 |
| 4 | 1100 |
| 5 | 1101 |
| 6 | 1110 |
| 7 | 1111 |

# Biased code - type 1

$$\boxed{\mathcal{A}_1(X) = X + \tfrac{1}{2}\mathcal{Z} - 1}$$

$$-\tfrac{1}{2}\mathcal{Z} < X \leq \tfrac{1}{2}\mathcal{Z} \qquad \text{☞} \qquad \textbf{asymmetric range}$$

$$\boxed{\begin{aligned} \textbf{MSB} = \textbf{1} &\iff X > 0 \\ \textbf{MSB} = \textbf{0} &\iff X \leq 0 \end{aligned}}$$

**z=2:**



| $X$ | $\mathcal{A}(1)\mathbf{X}$ |
|:---:|:---:|
| -7 | 0000 |
| -6 | 0001 |
| -5 | 0010 |
| -4 | 0011 |
| -3 | 0100 |
| -2 | 0101 |
| -1 | 0110 |
| 0 | 0111 |
| 1 | 1000 |
| 2 | 1001 |
| 3 | 1010 |
| 4 | 1011 |
| 5 | 1100 |
| 6 | 1101 |
| 7 | 1110 |
| 8 | 1111 |

**addition and subtraction** :

$$\mathcal{A}(A + B) = \mathcal{A}(A) + \mathcal{A}(B) - K$$
$$\mathcal{A}(A - B) = \mathcal{A}(A) - \mathcal{A}(B) + K$$

$$\mathcal{A}(A + B) = \mathcal{A}(A) + \mathcal{A}(B) - \tfrac{1}{2}\mathcal{Z} + 1$$
$$\mathcal{A}(A - B) = \mathcal{A}(A) - \mathcal{A}(B) + \tfrac{1}{2}\mathcal{Z} - 1$$

$$-\tfrac{1}{2}\mathcal{Z} \equiv +\tfrac{1}{2}\mathcal{Z} \quad (\textbf{mod } \mathcal{Z}) \qquad \tfrac{1}{2}\mathcal{Z} \sim \boxed{100\ldots00}$$

$-\tfrac{1}{2}\mathcal{Z} \equiv +\tfrac{1}{2}\mathcal{Z}$ **(mod $\mathcal{Z}$) $\rightsquigarrow$ negation of bit in higher order**

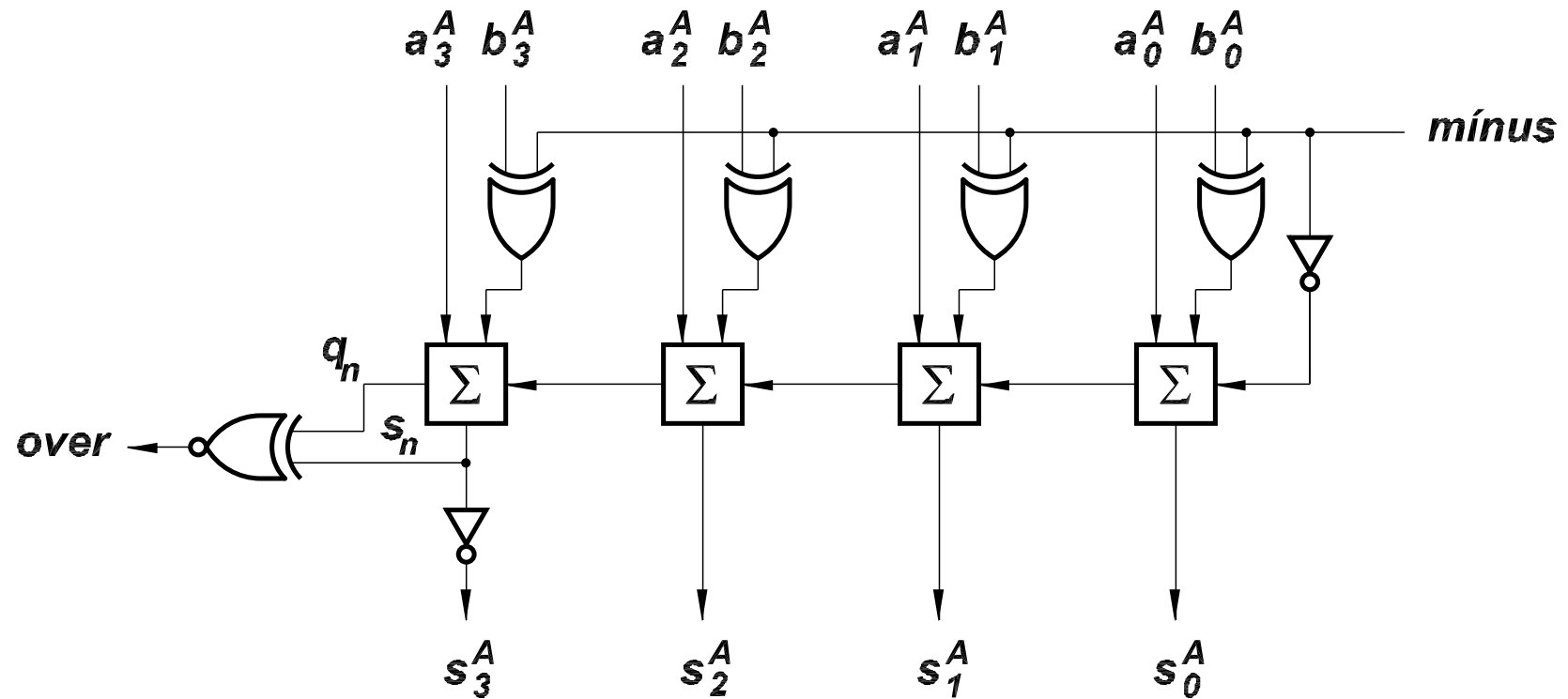*overflow:*

**addition:**     **same sign of addend**
          **and different sign of result**

**subtraction:   sign of minuend and subtrahend are different**
          **signs of minuend and result are different**

(derivation is analogous as in 2's complement code)

**z=2:**

**adder / subtractor**



**sign change** :      $-X$   =    $0 - X$

**overflow detection**

| $a_n$ | $b_n$ | $p_n$ | $q_n$ | $s_n$ | |
|:---:|:---:|:---:|:---:|:---:|:---|
| $\boxed{0}$ | $\boxed{0}$ | 0 | 0 | $\boxed{0}$ | $q_n = s_n$ |
| 0 | 0 | 1 | 0 | 1 | $q_n \neq s_n$ |
| 0 | 1 | 0 | 0 | 1 | $q_n \neq s_n$ |
| 0 | 1 | 1 | 1 | 0 | $q_n \neq s_n$ |
| 1 | 0 | 0 | 0 | 1 | $q_n \neq s_n$ |
| 1 | 0 | 1 | 1 | 0 | $q_n \neq s_n$ |
| 1 | 1 | 0 | 1 | 0 | $q_n \neq s_n$ |
| $\boxed{1}$ | $\boxed{1}$ | 1 | 1 | $\boxed{1}$ | $q_n = s_n$ |

$$over = \overline{q_n \oplus s_n}$$