# Information technology – lecture 10
## Number systems and representations. Computer arithmetic.

Roman Putanowicz
R.Putanowicz@L5.pk.edu.pl

# Positional notation

Let:
$\beta \in N, \beta >= 2$ – the base
$x_k$ – digits, $0 <= x_k < \beta$ with $k = -m, \ldots, n$

Notation:

$$x_\beta = (-1)^s [x_n x_{n-1} \ldots x_1 x_0 . x_{-1} x_{-2} \ldots x_{-m}] \quad x_n \neq 0$$

Interpretation:

$$x_\beta = (-1)^s \left( \sum_{k=-m}^{k=n} x_k \beta^k \right)$$

# Most common positional systems

decimal – base $\beta = 10$,
        digits: 0,1,2,3,4,5,6,7,8,9

binary – base $\beta = 2$,
        digits: 0,1

octal – base $\beta = 8$,
        digits: 0,1,2,3,4,5,6,7

hexadecimal – base $\beta = 16$,
        digits: 0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F

# Numbers in decimal system

$$175_{(10)} \rightarrow \text{digits} \rightarrow \quad 1 \mid 7 \mid 5$$

$$175_{(10)} = (-1)^0 \left( [\mathbf{1} \times 10^2] + [\mathbf{7} \times 10^1] + [\mathbf{5} \times 10^0] \right) \qquad (1)$$

# Numbers in binary system

$$175_{(10)} \rightarrow \text{digits} \rightarrow \quad 1 \mid 0 \mid 1 \mid 0 \mid 1 \mid 1 \mid 1 \mid 1$$

$$175_{(10)} = (-1)^0 \left( [\mathbf{1} \times 2^7] + [\mathbf{0} \times 2^6] + [\mathbf{1} \times 2^5] + [\mathbf{0} \times 2^4] + \right.$$
$$\left. [\mathbf{1} \times 2^3] + [\mathbf{1} \times 2^2] + [\mathbf{1} \times 2^1] + [\mathbf{1} \times 2^0] \right)$$

# Numbers in octal system

$$175_{(10)} \rightarrow \text{digits} \rightarrow \quad 2 \mid 5 \mid 7$$

$$175_{(10)} = (-1)^0 \left( [\mathbf{2} \times 8^2] + [\mathbf{5} \times 8^1] + [\mathbf{5} \times 8^0] \right)$$

# Numbers in hexadecimal system

$$175_{(10)} \rightarrow \text{digits} \rightarrow \quad \text{A} \mid \text{F}$$

$$175_{(10)} = (-1)^0 \left( [\mathbf{A} \times 16^1] + [\mathbf{F} \times 10^0] \right)$$

# Conversion binary $\rightarrow$ octal $\rightarrow$ decimal

| binary | 10101111 | | |
|---|---|---|---|
| | 10 | 101 | 111 |
| | 2 | 4+1 | 4+2+1 |
| octal | 257 | | |
| | $2 \times 8^2 + 5 \times 8^1 + 7 \times 8^0$ | | |
| decimal | 175 | | |

# Conversion decimal → other bases

While the quotient is not 0
    Divide the decimal number by the new base.
    Make the reminder the next digit to the left in the answer.
    Replace the decimal number with the quotient.

# Conversion decimal → binary

| 175 | integer division by 2: quotient, reminder | |
|---|---|---|
| 87 | 1 | rightmost digit |
| 83 | 1 | |
| 21 | 1 | |
| 10 | 1 | |
| 5 | 0 | |
| 2 | 1 | |
| 1 | 0 | |
| 0 | 1 | leftmost digit |
| 10101111 | | ← result |

# Conversion decimal → hexadecimal

| 175 | | | integer division by 16: quotient, reminder |
|-----|-----|-----|-----|
| | 10 | F | rightmost digit |
| | 0 | A | leftmost digit |
| AF | | | ← result |

# Arithmetic in binary system

## Addition

```
   84          01010100
 + 91        + 01011011
 ---          ---------
  175          10101111
```

## Multiplication

```
   35          00100011
 *  5        * 00000101
 ---          ---------
  175          00100011
               00000000
             + 00100011
               ---------
               10101111
```

# Fixed point representation of real numbers

Let us assume that a real number is represented with N memory positions such that 1 position is hold for sign, $N - k - 1$ positions for integer digits, $k$ positions for the digits after the point.

**Notation:**

$$x = (-1)^s \left[ a_{N-2} a_{N-3} \ldots a_k \cdot a_{k-1} \ldots a_0 \right]$$

**Interpretation:**

$$x = (-1)^s \beta^{-k} \sum_{j=0}^{N-2} a_j \beta^j$$

# Floating point representation of real numbers

Let us assume that a real number is represented with N memory positions such that 1 position is hold for sign, $N - k - 1$ positions for integer digits, $k$ positions for the digits after the point.

**Notation:**

$$x = (-1)^s \left[ 0 . a_1 a_2 \ldots a_t \right] \beta^e$$

**Interpretation:**

$$x = (-1)^s \times m \times \beta^{e-t}$$

where: t – the number of allowed significant digits $a_j$,

      m – integer number called mantisa

      e – integer number called exponent

The number zero has a separate representation.

# Floating point number formats

On 32 bit machine:

## Single precision

| 1 | 8 bits | 32 bits |
|---|---|---|
| s | e | n |

## Double precision

| 1 | 11 bits | 52 bits |
|---|---|---|
| s | e | n |

# Floating point number sets

Let us define a set of floating point numbers:

$$\mathcal{F}(\beta, t, L, U) = \{0\} \cup \{x \in R : x = (1)^s \beta^e \sum_{i=1}^{t} a_i \beta^{-}i\}$$

where: $t$ – number of significant digits
$\beta$ – base,
$0 < a_i \leqslant \beta - 1$ – digits
range $(L, U)$ such that $L \leqslant e \leqslant U$

# Normalisation of floating point number representation

To enforce uniqueness in number representation it is assumed that:

$$a_1 \neq 0$$
$$m \geqslant \beta^{t-1}$$

Such representation is called **normalized**.
With such setup $a_1$ is the primary significant digit, $a_t$ last significant digit.

# Distribution of floating point numbers

Floating point numbers are not equally spaced along the real line

machine epsilon – the smallest number in the set of floating point numbers such that $1 + \varepsilon_m > 1$

# IEEE 754 Standard

By choosing different set of values for $\beta, t, L$ and $U$ it is possible to build multiple floating point number systems. Thus the necessity of defining a standard for floating point arithmetic. One of the most widely-used is IEEE Standard for Floating Point Arithmetics (IEEE 754). The standard defines:

- arithmetic formats (binary and decimal)
- interchange formats,
- rounding algorithms,
- operations (arithmetic and other)
- exception handling