

Lista 3

tags: ASK

#Zadanie 2

Zmiennopozycyjna reprezentacja liczby: $(-1)^S 2^E M$, gdzie $E = \text{EXP} - \text{bias}$.

16 bitowe liczby zmiennopozycyjne w standardzie IEEE 754-2008:

- $S < 1$ bit
- $\text{EXP} < 5$
- $M < 10$ (wiodąca jedynka jest w domyśle)
- $\text{bias} = 15$

$$0.15625_{(10)} = 0.00101_{(2)} = 1.01 \cdot 2^{-3}$$

$$(\text{EXP} = E + \text{bias} = -3 + 15 = 12 = 1100_2)$$

Kodując dostajemy:

Liczba	S	EXP	M
$1.01 \cdot 2^{-3}$	0	01100	0100000000

Półowiczna precyzja (16 bit)

- S_{\min} : 0 00001 0000000000
- S_{\max} : 0 11110 1111111111

Zakres: $[6.1 \cdot 10^{-5}, 65504]$ (jeśli zdenormalizowane to minimum wynosi $5.96 \cdot 10^{-8}$)

Pojedyncza precyzja (32 bit)

- S_{\min} : 0 00000001 0000000000000000000000
- S_{\max} : 0 00000000 111111111111111111111111

Zakres: $[1.7510^{-38}, 3.410^{38}]$ (jeśli zdenormalizowane to minimum wynosi $1.4 \cdot 10^{-45}$)

Zadanie 3

Reprezentacja mantysy w terminach S_{guard} , S_{round} , S_{sticky}

$M = 1.\text{BBB}..\text{BGRXXX}..X$ liczba bitów $B + 1$ to rozmiar mantysy, gdzie $G(\text{guard})$ to ostatni bit, który zachowujemy przy zaokrągłaniu. $R(\text{round})$ to pierwszy tracony bit (połówka). $X(\text{sticky})$ to każdy kolejny bit po R .

(bias = 15)

Liczba	S	EXP	M	E
0.3984375	0	01101	1001100000	-2
0.34375	0	01101	0110000000	-2
1771	0	11001	1011101011	10

Dodając od lewej $0.3984375 + 0.34375 = 1.001100000 \cdot 2^{-2} + 1.011000000 \cdot 2^{-2} = 10.1111100000 \cdot 2^{-2} = 1.01111100000 \cdot 2^{-1}$

$$+1771 \cdot 1.01111100000 \cdot 2^{-1} + 1.1011101011 \cdot 2^{10} = 0.000000000010111112_{(10)} + 1.1011101011 \cdot 2^{10} = 1.1011101011101111 \cdot 2^{10}$$

S_{round} i S_{sticky} to 1 , więc zaokrąglamy w górę $1.1011101011 \cdot 2^{10} \approx_r 1.1011101102_{(10)} = 1772_{(10)}$

Dodając od prawej $0.34375 + 1771 = 1.011000000 \cdot 2^{-2} + 1.1011101011 \cdot 2^{10} = 0.00000000001011 \cdot 2^{10} + 1.1011101011 \cdot 2^{10} = 1.101110101101011 \cdot 2^{10} \approx_r 1771_{(10)}$

$$+0.3984375 = 0.000000000001001 \cdot 2^{10} + 1.1011101011 \cdot 2^{10} = 1.101110101101001 \cdot 2^{10} \approx_r 1771_{(10)}$$

Zadanie 4

```
int32_t x, y, z;  
dx = (double)x;  
dy = (double)y;  
dz = (double)z;
```

```
//Double dokładnie reprezentuje każdą zmienną typu int32_t
```

- $(float)x == (float)dx$ Wyrażenie zawsze oblicza się do prawdy. Castowanie do $float$ będzie się wiązać z taką samą utratą bitów.
- $dx - dy == (double)(x - y)$ W operacji $x - y$ może wystąpić zjawisko $overflow$ jeśli odejmujemy liczby z zakresu `int32_t` (np. $x = INT32_MIN$, $y = 1$). Natomiast dla $dx - dy$ wartość wyliczy się poprawnie.
- $(dx + dy) + dz == dx + (dy + dz)$ Wyrażenie zawsze oblicza się do prawdy. Wyniki tych działań zmieszczą się w reprezentacji `double`.
- $(dx * dy) * dz == dx * (dy * dz)$ Wyrażenie może obliczyć się do fałszu, ponieważ nie jesteśmy w stanie reprezentować dokładnie wyników mnożenia dowolnych dwóch liczb z zakresu typu `int`, czyli będziemy wykonywać zaokrąglenia. Przykładowo dla wartości $x = 6$, $y = 2^{23} + 3 + 1$, $z = 2^{30}$.

Od lewej `` 6 x 8 388 611

50 331 666 binarnie: 11000000000000000000010010 <- mieści się w mantysie

50 331 666 x 1 073 741 825 = 54 043 214 906 130 450 binarnie: 110000000000000000000100100000110000000000000000010010

Od prawej

1 073 741 825 x 8 388 611 = 9 007 202 484 355 075 binarnie:

100000000000000000000000011000000100000000000000000000000000011

Zaokrąglamy do parzystej:

10000000000000000000000001100000010000000000000000000000000010

1000000000000000000000000110000001000000000000000000000000010_2

x 110_2

110000000000000000000000010010000011000000000000000000001100

- `***dx / dx == dz / dz***`

Dla $x = 0$, $y = 1$ lewa strona oblicza się do NaN , a prawa nie.