

LISTA 2

~~1, 2, 3, 4, 5, 6, 7, 8~~

2000.00

34

2AD. 1
1) $(x > 0) \parallel (x - 1 < 0)$

Nie, ponieważ dla $x = -2$ prawa i lewa strona wyrażenia jest fałszywa.

$x = -2$ nie jest większy od 0, a
 $x = -1 = -2$ nie jest mniejszy od 0.

0. 1 1 1 1 1 1 1
 100000 . 000 ←
- 000000 . 001

 011111 111 } max
 } dad

$$\begin{aligned} 0-0 &= 0 \\ 0-1 &= 1+1 \\ 1-0 &= 1 \\ 1-1 &= 0 \end{aligned}$$

2) $(x \neq 7) \equiv 7 \vee (x < 29 \wedge 0)$ $7 = 111$

Tak, 7 w zapisie binarnym to 111, zatem jedyną sytuację w której prawa strona wyrażenia jest fałszywa to taka gdy na ostatnich trzech pozycjach x będzie miały same jedynki. Przedstawiając takiego x pod prawą stronę wyrażenia zawsze będzie ono prawdziwe, gdyż po przemieszczeniu x będzie miało postać 110...0, a więc będzie liczbą parzystą.

$$3) (x * x) \geq 0$$

$$4) (x < 0) \parallel (-x \leq 0)$$

- Tak, rozważmy dwa przypadki:
- gdy $x < 0$ to spełniona jest lewa część wyrażenia
 - gdy $x \geq 0$ to spełniona jest prawa część wyrażenia gdyż $|\text{int32_t_MAX}| < |\text{int32_t_MIN}|$ więc nie przekracza zakresu.

$$5) x > 0 \parallel -x \geq 0 \quad \text{Nie, dla } x = \text{INT32_MIN}$$

- Nie, rozważmy dwa przypadki:
- gdy $x > 0$ to spełniona jest lewa część wyrażenia
 - gdy $x \leq 0$ to zrobi fałsz, ponieważ $|\text{int32_t_MIN}| > |\text{int32_t_MAX}|$, czyli przekrocza zakres.

$$6) (x \mid -x) \gg 31 == -1$$

Nie, dla $x = 0$, bo $0 \gg 31 = 0$
(dla wszystkich innych jest spełniona)

$$7) ((\text{uint32_t})x) \gg 3 == ((\text{uint32_t})x) / 8$$

$$2^3 = 8$$

Tak, ponieważ przesunięcie bitowe to to samo co dzielenie, dla liczb dodatnich (bez znaku)

$$8) x \gg 3 == x / 8 \quad x = -1 \quad \text{jak powiedzieliśmy to powyżej to 0, a przesunięcie nie robi}$$

W naszym przypadku
zajmowane bity

Nie, ponieważ dla liczb ujemnych (ze znakiem) przesunięcie w prawo spowoduje skopiowanie najbardziej znaczącego bitu na nową pozycję. Wynikiem będzie zawsze fałsz dla ujemnych liczb nieparzystych (zaczynając od pierwszego bitu)

$$9) x \gg 2 == (\text{uint32_t})x + (\text{uint32_t})x$$

Tak, podczas wykonywania operacji zmienną ze znakiem zostają zmiennione bez znaku.

2AD.2

Prosty sposób na sprawdzenie czy $x < y$ to obliczenie $x - y$. Gdy ta wartość jest ujemna, to y jest większy, wpp. nie jest większy.

Niestety w ten sposób możemy uzyskać overflow np. $\text{INT_MIN} - 1$ da nam liczbę dodatnią, a w tym przypadku $x < y$, więc mamy błąd. Rozwiązaniem jest podzielenie wartości x i y przez 2.

To tworzy kolejny problem, ponieważ tracimy bity jedynki w x i y .

Przykład:

A) x nie zawiera 1 i y nie zawiera 1:
nic nie robimy, ponieważ nic nie tracimy

B) x nie zawiera 1 i y zawiera 1: odejmujemy od notacji ($\sim x \& y \& 1$), ponieważ w tej sytuacji y mogło być większe od x , ale straciliśmy bit który o tym przesądzał. Dokładniej $(x > 1) - (y > 1)$ będzie równe 0, a my odejmujemy od tego 1 i dostaniemy -1

C) x zawiera 1 i y nie zawiera 1: nic nie robimy, ponieważ wtedy jeśli $(x > 1) - (y > 1)$ będzie równe 0, to i tak $x > y$

D) x zawiera 1 i y zawiera 1: nic nie robimy, ponieważ wtedy jeśli $(x > 1) - (y > 1)$ będzie równe 0, to $x = y$

Następnie przechwycamy bit znaku o 31 i jeśli mamy ujemną to w operacji $\&$ de

$$((x > 1) - (y > 1) - (\sim x \& y \& 1)) > 31) \& 1$$

2AD.3

$b = x > 31$
 $\text{return } \underbrace{\sim b \& x}_{\text{dodat. } b = 000\dots 0} + \underbrace{b \& \sim(x-1)}_{\text{ujemne. } b = 111\dots 1}$

<p>dodat. $b = 000\dots 0$</p> <p>ujemne. $b = 111\dots 1$</p> <p>zero $b = 00\dots 0$</p>	<p>x</p> <p>$0000\dots 0$</p> <p>x</p>	<p>$+$</p> <p>$00000\dots 0$</p> <p>$-x$</p> <p>$00\dots 0$</p>	<p>$= x$</p> <p>$= -x$</p> <p>$= x$</p>
---	---	---	--

$\sim(x-1) = x \text{ np. } x = 1001 (-7) = 0001 = 1000 = 0111 \rightarrow 7$

1NN4

$$y = x > 31$$

$$\text{return } (x \wedge y) - y$$

2AD.4 $(x > 31) - ((-x) > 31)$

dla $x < 0$

$$\underbrace{1\dots 1111} - \underbrace{000\dots 0} = \underbrace{11111111}_{\text{jeśli się dowiódł, że min się nie zmieści to napisz min32 + y = -x i podstaw y podl - x}} = -1$$

dla $x = 0$

$$00000\dots 0 - 000\dots 0 = 000\dots 0 = 0$$

dla $x > 0$

$$000\dots 0 - 1\dots 1111111 = 000\dots 01 = 1$$

$$\begin{array}{r} 111 \\ 0000 \\ 1111 \\ \hline 0001 \end{array}$$

ZAD 5

```
int32_t z5_inrange(int32_t x, int32_t y)
{
    int32_t z = x + y;

    int32_t a = x >> 31;
    int32_t b = y >> 31;
    int32_t c = 2 >> 31;

    return ((a ^ c) & (b ^ c)) + 1;
}
```

3 przypadki:

- x, y i z mają taki sam znak wtedy nie ma niedomiaru ani nadmiaru, więc zwraca 1
- x, y mają ten sam znak, a z ma inny wtedy mamy nadmiar lub niedomiar
- jeśli x i y mają różne znaki to nie ma nadmiaru ani niedomiaru

ZAD. 6

Na podstawie zad. 4 z listy 1

```
int32_t odd_one(int32_t x) {
```

```
    x = (x & 0x55555555) + ((x >> 1) & 0x55555555);
    x = (x & 0x33333333) + ((x >> 2) & 0x33333333);
    x = (x & 0x0f0f0f0f) + ((x >> 4) & 0x0f0f0f0f);
    x = (x & 0x00ff00ff) + ((x >> 8) & 0x00ff00ff);
    x = (x & 0x0000ffff) + ((x >> 16) & 0x0000ffff);
    return x & 0x00000001;
}
```

}

Zad. 7 zwraca ^{bit} bitów które mają ile 1 występuje w x.
Ostatnie operacje sprawdza co znajduje się na ostatnim bicie tej liczby. Gdyż to mówi nam czy liczba jest parzysta czy nieparzysta.

ZAD. 8

1) $x == (\text{int32_t})(\text{double})x$

1) TAK

double ma taką samą precyzję jeśli chodzi o liczby całkowite co

int32_t \rightarrow nie większą precyzję

2) NIE, float ma 32 bity, ale musi mieć więcej cyfr całkowitych i ułamkowych

3) $d == (\text{double})(\text{float})d$

float ma mniejszą liczbę bitów, więc ma mniejszy zakres

4) $f == (\text{float})(\text{double})f$ ❌

Proceda

5) $f == -(-f)$

Proceda. float w swojej reprezentacji trzyma minus. Nie najstarszym bicie. Znak "-" mówi nam znak na najstarszym bicie

6) $1.0 / 2 == 1 / 2.0$

Proceda, jeżeli chociaż jedno float to drugie też zostanie przerobione na float

7) $d * d \geq 0.0$

Proceda

Jeżeli f jest bardzo duże, o d może to f+d przekroczyć się do f i wtedy $f-f=0$

8) $(f + d) - f == d$

a d nie jest 0

Nie, ze względu na niedokładność obliczeń ~~nie~~ na werbach zmiennopozycyjnych

ZAD. 7

Standard IEEE 754-2008 definiuje liczbę zmiennopozycyjną o szerokości 16-bitów. Oblicz ręcznie $3.984375 \cdot 10^{-1} + 3.4375 \cdot 10^{-1} + 1.741 \cdot 10^3$ mając liczbę w tym formacie. Zapisz wynik binarnie i dziesiętnie. Czy wynik się mieści? Jeśli najpierw wykonamy drugie obliczenie?

ROZWIĄZANIE

Liczba zmiennopozycyjna 16-bitowa w standardzie IEEE 754-2008 składa się z 1 bitu znaku, 5 bitowego wykładnika oraz 10 bitowej mantysy. Bicie wykładnika wynosi 15.

1. KONWERSJA

Zaczniemy od zamiany liczb do postaci właściwej dla rozważanego sposobu zapisu liczb zmiennopozycyjnych.

$$3.984375_{10} \cdot 10^{-1} = 0.3984375_{10} = 0.0110011_2 = 1.10011_2 \cdot 2^{-1}$$

$$\frac{1}{2} = 0.5 \quad \frac{1}{4} = 0.25 \quad \frac{1}{8} = 0.125 \quad \frac{1}{16} = 0.0625$$

$$0.3984375 = 0 \cdot (1/2) + 1 \cdot (1/4) + 1 \cdot (1/8) + \dots$$

To liczba w formacie 16-bitowym IEEE 754-2008 zostałaby zapisana następująco:

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
+	-	WYKŁADNIK					MANTYSA								
0	0	1	1	0	1	1	0	0	1	1	0	0	0	0	0

Podobnie postępujemy z pozostałymi liczbami
 $3.4375_{10} \cdot 10^{-1} = 0,34375_{10} = 0,01011_2 = 1,011_2 \cdot 2^{-2}$
 $1.471_{10} \cdot 10^3 = 11011101011_2 = 1,1011101011_2 \cdot 2^{10}$

Należy sprawdzić czy:

- część ułamkowa mieści się w 10-bitowej mantysie (tak, w szczególności ostatnia liczba w części ułamkowej mantysy)
- wykładnik zawiera się w przedziale -14..15 (tak)

2. DODAWANIE

Dodawanie liczb zmiennopunktowych odbywa się przez zbadanie pozycji separatora części całkowitej i ułamkowej, wykonanie działania, bez straty (dodatkowo) i ostatecznie konwersji do używanego formatu.

$$(1,10011_2 \cdot 2^{-2} + 1,011_2 \cdot 2^{-2}) + 11011101011_2 \cdot 2^{10} =$$

$$= 1,011111_2 \cdot 2^{-1} + 1,1011101011_2 \cdot 2^{10} =$$

$$= 11011101011,0110011_2$$

$$\begin{array}{r} 0 \quad 0110011 \\ + 0 \quad 0101100 \\ \hline = 0 \quad 1011111 \end{array}$$

$$\begin{array}{r} 11011101011 \quad 0000000 \\ + 00000000000 \quad 1011100 \\ \hline = 11011101011 \quad 1011100 \\ \quad \quad \quad G \quad R \quad X \quad X \quad X \quad X \quad X \end{array}$$

Mantysa mieści jednak tylko 10 bitów - musimy dokonać zaokrąglenia. Bity G (ostatni w mantysie) wynosi 1, R (pierwszy po zero) wynosi 1, zatem zaokrąglenie wykonamy w górę.

Gdyby mieliśmy kolejność dodawania $1,10011_2 \cdot 2^{-2} + (1,011_2 \cdot 2^{-2} + 1,1011101011_2 \cdot 2^{10})$ pierwsze dodawanie przebiegłoby by następująco

$$\begin{array}{r} 11011101011 \quad 0000000 \\ + 00000000000 \quad 0101100 \\ \hline 11011101011 \quad 0101100 \\ \quad \quad \quad G \quad R \quad X \quad X \quad X \quad X \quad X \end{array}$$

w tej sytuacji G=1, R=0, S=1 - wykonujemy zaokrąglenie w dół