

**From:** [Lukasz Tomaszek](#)  
**To:** [psokolnicki@interia.pl](mailto:psokolnicki@interia.pl)  
**CC:**  
**Subject:** iee 754  
**Date:** 25 stycze• 2006 17:08:13  
**Attachments:** [ATT00005.txt](#)

---

## Co student PWR o IEEE 754 wiedziec powinien

### 0. Wstep

---

To bedzie szybka manual na temat IEEE 754, dzialan i nawiazan do rzeczywistosci (systemu dziesietnego). Nikomu nie bede tlumaczyl zapisu liczby w U2, bez znaku czy innych, prozaicznych rzeczy. Rozsadne pytania i sugestie mile widziane. Jesli nie rozumiesz prostych zapisow, to dla Ciebie juz nie ma nadziei.

Tekst pisany jest 'bez polskich znaczkow'. Byc moze kiedys sie to zmieni.

#### ChangeLog:

18/19 marca 2003 - wersja 0.8 RC1

Pierwsza poprawna wersja publiczna. Tekst nie jest jeszcze skonczony. Wersja 'przymusowa' zawierala wiele bledow i sprzecznosci, ktore staralem sie poprawic.

14 marca 2003 - wersja 0.2B

Pierwsza (przymusowa) wersja publiczna. Strescilem rozdzialy 2 i 3. Sypnal mi sie komp i rozdzialy 4 i kawalek 3 i 5

12 marca 2003 - wersja 0.1B

Pierwsza wersja. Jedyne rozdzialy 0 i 1.

#### ToDo:

- Dodawanie i odejmowanie (w razie koniecznosci **czytajcie archiwum InfEki**)
- Zaokraglanie do nieparzystej
- Obliczanie bledu sredniego zaokraglenia
- Polskie znaki w tekście i poprawa bledow ortograficznych/gramatycznych
- Uzupełnienia

© IronHand & Ryan O'Connell

### 1. Liczba

---

Liczba zmiennoprzecinkowa pojedynczej precyzji (**single** - nie mylic z liczba podwojonej precyzji - double, o tych samych, zreszta, czesciach) sklada sie ze **znaku**, **mantysy** i **wykladnika**. Znak przyjmuje wartosc **0** dla liczby dodatniej i **1** dla liczby ujemnej. Znak oznaczac bedziemy **Z**. Mantysa to cyfry po przecinku, przed ktorymi domyslnie stoi jedynka (nie zawsze - o tym dalej). Mantysy zapisujemy na **23** bitach. Przyklad mantysy i co ona oznacza:

```
      / ----- 23 bity ----- \  
Mantysa 1010101010101010101010101  
Oznacza 1,1010101010101010101010101  
  
Mantysa 00000000000000000000000  
Oznacza 1,00000000000000000000000
```

Warto zaznaczyc, ze jest to liczba **bez znaku** o ewentualnej ujemnosci informuje nas znak kodowany osobno. Dlatego tez same zera w mantysie to zera po przecinku i **1** przed przecinkiem, czyli w systemie dziesietnym jeden (ta jedynka na poczatku nie oznacza liczby ujemnej jak w U2; jesli to Cie konfuduje - dopisz 0 na poczatek). Mantysy oznaczali bedziemy **M**. Ostatnia 'czescia skladowa' liczby zmiennoprzecinkowej jest wykladnik. Oznaczac bedziemy go **W**. Wykladnik zapisywany jest na **8** bitach w systemie **+127** (z obciazeniem 127). Oznacza to, ze po odczytaniu wartosci liczby zapisanej na 8 bitach jako liczby bez znaku by uzyskac wartosc wykladnika trzeba odjac 127 (bo liczba jest obciazona o wartosc 127). Przyklady:

```
Wykladnik 00000010  
Dziesietnie      2  
Wartosc W  2 - 127 = -125  
  
Wykladnik 01111111
```

```

Dziesietnie      127
Wartosc W      127 - 127 = 0

Wykladnik      10000111
Dziesietnie      135
Wartosc W      135 - 127 = 8

```

Czyli wykladnik moze osiagac wartosci od -127 do +128.

Aby poznac wartosc liczby w IEEE754 nalezy wszystko to zebrac do kupy. Oznacza to, ze trzeba uwzglednic znak a mantyse pomnozyc przez 2 do potegi wykladnik. Przyklad:

```

zawartosc      oznacza
Z 0            +1
M 101010101010101010101010101010101  1,101010101010101010101010101010101
W 011111111    0
L = +1,101010101010101010101010101010101 * 20 = +1,101010101010101010101010101010101

```

```

zawartosc      oznacza
Z 0            +1
M 000000000000000000000000000000001  1,000000000000000000000000000000001
W 10000001     2
L = +1,000000000000000000000000000000001 * 22

```

```

zawartosc      oznacza
Z 1            -1
M 000000000000000000000000000000000  1,000000000000000000000000000000000
W 00001000    -119
L = -1,000000000000000000000000000000000 * 2-119 = -1 * 2-119

```

Oczywiscie w prawdziwych maszynach nie przechowujemy wszystkich tych wartosci osobno. Liczba bitow, na ktorych mozna przechowac liczbe zmiennoprzecinkowa wg IEEE754 to 32 (bo 1+8+23). Kolejnosć zapisu wartosci takze jest ustalona i stanowi **ZWM** (ZWM jak Zwiasek Walki Mlodych - taki twor; latwo zapamietac). Przykladowe liczby:

```

      1   1
( 1 + - + - ) * 2
      2   8
|Z|      W      |      M      |
0 1 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
      /  |  \
      /  |  \
      1/2 1/4 1/8 itd. te miejsca maja wage:
te miejsca maja wage:
1/2 1/4 1/8 itd. Dlaczego? Architektura Komputerow sie klania

-256 (czyli -1*28)
|Z|      W      |      M      |
1 1 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

```

Wspomnialem wczesniej o liczbach, ktore nie maja domyslnej 1 przed mantysa. Sa to tak zwane liczby zdenormalizowane (w przeciwienstwie do tych, ktore jedynke maja i mienia sie znormalizowanymi). Oprócz podzialu na liczby z(de)normalizowane mamy takze kilka wartosci specjalnych. Wartosci te to 0 (a raczej +0 i -0), + i - nieskonczonosc (bedziemy je oznaczac +inf i -inf) oraz nieliczba (wartosc, ktora nie podpada pod nic powyzzszego, np. wynik dzielenia 0/0) oznaczana NaN. Wszystkie powyzzsze wartosci poznajemy po wartosciach mantysy i wykladnika. Obrazuje to ponizsza tabela:

W	M	Typ liczby
rozny od 0 i rozny od MAX	dowolna	Liczba znormalizowana
0	rozna od 0	Liczba zdenormalizowana
0	0	+/-0 (znak zalezny od bitu znaku)
MAX	rozna od 0	NaN (rozroznianie znaku nie ma sensu)
MAX	0	+/-inf

Pod pojeciem MAX rozumie same 1 (jedynki na wszystkich polach M lub W zaleznie od przypadku).



```
Mantysa      1,11100001111000011110000 przesuwamy o cztery pozycje,
Po przesunieciu 0,00011110000111100001111000
```

Teraz wracamy do zaokrąglenia. Zaokrąglenie przez obcięcie polega na obcięciu części, która nie mieści się na pozycji mantysy. Nie warto dalej tłumaczyć, gdyż nie ma czego. Przykład:

```
          / ----- 23 bity ----- \
Wartosc do zaokrąglenia 0,10000000000000000000000100
Wartosc po zaokrągleniu 0,10000000000000000000000000

Wartosc do zaokrąglenia 0,00001010101010101010101011
Wartosc po zaokrągleniu 0,0000101010101010101010101
```

Krotko mowiac: nie wazne co bylo w polach GRS - odcinamy to. Taki sposob zaokrąglania nie jest jednak stosowany kiedykolwiek (choc byloby prawdopodobniej nam latwiej). Odciecie nie jest w ogole czescia standardu, jednak stare układy arytmetyczne, które proklamowały swa kompatybilnosc ze standardem IEEE 754 zanim zostal on jeszcze zatwierdzony, implementowały ja jako glowna (jesli nie jedyna) metode zaokrąglania. Takie 'cacka' mozna nazwac starymi procesorami arytmetyki zmiennoprzecinkowej.

### Zaokrąglenie do +inf

To zaokrąglenie przyjmuje rozna postac zaleznie od znaku liczby. Latwo to wytłumaczyc tym, ze liczba ujemna zbliżając się w strone nieskonczonosci maleje w module (maleje jej wartosc z pominięciem znaku), podczas gdy liczba dodatnia rośnie. Ilustruje to ponizsza zalezność. Przedstawilem ja wedlug zasady: liczbe z ułamkiem dziesietnym zaokrąglamy (tak jak w szkole, w zerowce) w gore (czyli w strone **+inf**):

```
Przed zaokrągleniem:
  -10,75      0      7,66
-----*-----|-----*----->

Po zaokrągleniu:
   -10      0      8
-----*-----|-----*----->
```

Jak widac sprawdzily sie przewidywania: liczba ujemna zmienila swoja wartosc modulo, liczba dodatnia zwiększyła. Podobnie dzieje sie przy zaokrąglaniu liczb zmiennoprzecinkowych zgodnych z IEEE 754.

**Jesli liczba jest ujemna to ucinamy jej nadmiarowa czesc** (podobnie jak odbywalo sie to w przypadku poprzednio omawianego zaokrąglania). Kiedy liczba jest dodatnia - zwiększamy mantysę o 1/2 ulp (czyli dodajemy jeden na pozycji mniejszej niz najmlodsza pozycja mantysy). Pozycja ta to **G**. Zaokrąglenie takie mozna stosowac jedynie gdy **G=1, R=0, S=0** (choc **nie** jest to zalecane) a nalezy wykorzystywac, gdy program wywołuje zaimplementowana funkcje **ceiling** (czyli zaokrąglanie liczby do najmniejszej liczby całkowitej nie mniejszej niz dana liczba). Przykładowo w przypadku **GRS = 100**, jesli dodamy jeden na pozycji **G** (gdzie juz jest **1**), to nastapi przeniesienie na pozycje starsza, czyli na pozycje najmlodsza mantysy (ta o wadze **ulp**). Ilustruje to przyklad:

```
          / ----- 23 bity ----- \
M 000000000000000000000001100
+1/2 ulp .....1..
-----
= 0000000000000000000000010000

          / ----- 23 bity ----- \
M 11111001111001110011010100
+1/2 ulp .....1..
-----
= 11111001111001110011011000
```

W powyższych przykladach przebieglo to bezstresowo. Co jednak, gdy mantysa zawiera same jedynki (a moze sie tak zdazyć)? Liczbe (po wykonaniu dzialania) nalezy przeskalowac, to znaczy tak zmienic jej mantysę i wykladnik, by przed mantysa znajdowala sie tylko domyslana jedynka. Pamietamy, ze przesuwajac mantysę o jeden bit w prawo nie zmieniamy wartosci liczby jesli inkrementujemy wartosc wykladnika i na odwrot (oczywiscie: nie zmieniamy wartosci, jesli jakis wazny, czyli zawierajacy jedynke, bit nie przeniesie nam sie na pola GRS). Ilustruje to przyklad:

```
          / ----- 23 bity ----- \
M 1,111111111111111111111111100
+1/2 ulp 0,00000000000000000000000100
-----
= 10,0000000000000000000000000000 przesuwamy o jeden bit w prawo,
  1,000000000000000000000000000000 nalezy jeszcze dodac jeden do wykladnika!
```

Tak oto nasza liczba ponownie stała się liczbą znormalizowaną.

Pozwól sobie jeszcze wrócić do tematu przeskalowywania. Otóż korzystając z kilku prostych operacji (mnożenia, dodawania i przesuwania bitów) można udowodnić, że:

```
M * 2W =  
M / 2 * 2 * 2W =  
M / 2 * 2W + 1 =  
M >> 1 * 2W + 1 =
```

Gdzie >> to przesunięcie bitów w prawo o pewną liczbę bitów/pozycji (tutaj: 1).

Oczywiście w powyższym przykładzie (i w kilku kolejnych) pisać **M** mam na myśli mantysę z dodaną już na początku jedynek (a nie goła sekwencję zer i jedynek wpisanych w pamięć). Dodatkowo tutaj zrezygnowałem z -127 w wykładniku w celu lepszego unaocznienia samego skalowania.

### Zaokrąglenie do -inf

Tutaj sytuacja jest podobna do zaokrąglania do **+inf**. Jedynie przypadki są 'odwrotne'. Otóż, jeśli zaokrąglimy liczbę ujemną w dół, to musimy dodać na pozycji G 1/2 **ulp** a jeśli zaokrąglimy liczbę dodatnią to stosujemy ucięcie.

Powód jest podobny do poprzedniego: liczba ujemna zaokrąglona w dół zwiększa swoją wartość w module, a liczba dodatnia zmniejsza:

```
Przed zaokrągleniem:  
-10,75            0            7,66  
-----*-----|-----*----->  
  
Po zaokrągleniu:  
-11               0            7  
-----*-----|-----*----->
```

Samo zaokrąglenie przebiega bardzo podobnie, więc ograniczę się do dwu przykładów:

```
Liczba dodatnia:  
  / ----- 23 bity ----- \  
  M 1,000000000001110000111111100  
Ucięcie 1,0000000000011100010000000  
  
Liczba ujemna (Z=1):  
  / ----- 23 bity ----- \  
  M 1,111111111111111111111111100  
+1/2 ulp 0,00000000000000000000100  
-----  
  = 10,00000000000000000000000000000000 przesuwamy o jeden bit w prawo,  
  1,00000000000000000000000000000000 należy jeszcze dodać jeden do wykładnika!
```

### Zaokrąglenie do zera

Z teoretycznego punktu widzenia ta metoda niczym nie różni się od zaokrąglania przez obcięcie. Dlaczego więc jest 'czymś innym'? Ponieważ generuje inne wyniki w systemach, w których podstawa nie jest **2**, lub znak zapisywany jest w alternatywny sposób (tzn. nie jako osobny bit). To jedyny powód rozróżniania tych dwóch nazw. W przypadku IEEE 754 z. do zera wygląda dokładnie tak, jak obcięcie. Kiedy stosujemy? Przy niektórych funkcjach konwersyjnych (np. przejście z float do integer, choć można stosować inne metody - wybór powinien należeć do użytkownika).

### Zaokrąglenie do najbliższej

Czym jest zaokrąglenie do najbliższej? Posłuży się analogia z systemu dziesiętnego. 10,2 jest bliżej 10 niż 11, więc zaokrąglamy do 10. 3,89 jest bliżej 4 i tak też zaokrąglamy. Jeśli ktoś jest bystry to zauważy, że wystarczy dodać połowę wartości na najmłodszej pozycji (zwaną ulp) i obciąć część ułamkową. I tak też działa to zaokrąglenie:

```
  / ----- 23 bity ----- \  
  M 1,10001000100010001000100011  
+1/2 ulp 0,00000000000000000000100  
-----  
  1,10001000100010001000100011 =  
  1,10001000100010001000100010
```

```

/ ----- 23 bity ----- \
M 1,111111111111100000000000101
+1/2 ulp 0,00000000000000000000100
-----
1,11111111111110000000001001 =
1,11111111111110000000001

```

Zaokrąglenie to stosowane jest jako domyślne z. w układach kompatybilnych z IEEE 754.

### Zaokrąglenie symetryczne do parzystej

Zaokrąglenie to nie znajduje się w specyfikacji IEEE 754, jednak implementowane jest często w układach arytmetyki zmiennoprzecinkowej. Jego drugą zaletą jest uwielbienie, jakim palają wykładowcy uniwersyteccy do tego rozwiązania. Do zalet z. symetrycznego do parzystej należy zerowy średni błąd zaokrąglenia (mówiąc prościej: jeśli rozpatrzemy wszystkie możliwe zaokrąglenia - czyli dla każdej kombinacji GRS - i obliczymy błąd, czyli różnicę między liczbą zaokrągloną a niezaokrągloną, to suma wszystkich błędów wyniesie zero; wynika to z faktu, że przy tym zaokrągleniu równie często występuje niedomiar, co i nadmiar, tzn. znak wyniku różnicy liczby przed i po zaokrągleniu będzie czasem + a czasem -).

Technika zaokrąglania tą metodą nie należy do prostych. W przypadku, gdy **RS** = 00 z. to wygląda bardzo prosto: jeśli **G** = 1 (a tylko wtedy można mówić o zaokrągleniu) dodajemy lub odejmujemy **1/2 ulp** zależnie od wartości najmłodszej pozycji mantysy (to znaczy jeśli mamy zero, a zero jest liczbą parzysta - odejmujemy, a gdy jedynek to dodajemy). Przykład wygląda tak:

```

/ ----- 23 bity ----- \
Wartosc do zaokrąglenia 0,1000000000000000000000100
Wartosc po zaokrągleniu 0,1000000000000000000000000
Wartosc do zaokrąglenia 0,10000000000000000000001100
Wartosc po zaokrągleniu 0,1000000000000000000000010

```

Zakładam, że każdy już wie jak przebiega dodawanie **1/2 ulp**.

Pozostałe przypadki nie są już takie przyjemne. Można tutaj wykorzystać jeden z dwóch schematów przedstawionych w Arytmetyce Komputerów znanego nam wszystkim autora. Przedstawię ten, który jest nieco łatwiej zapamiętać (bo oba są niezwykle podobne).

Wybieramy liczbę nie mniejszą niż zaokrąglana, która na najmłodszym bicie mantysy ma **1**. Liczba ta ma być możliwie zbliżona do liczby zaokrąglanej, np. jeśli zaokrąglamy liczbę kończąca się na **...0110** (trzy kolorowe bity to **GRS**) to bliska jej liczba jest **...1xxx** (**GRS** nas nie interesują - nie ma ich). Teraz obliczamy różnicę: liczba zaokrąglana - liczba przez nas wybrana. Tutaj różnica będzie **+0001** (gdzie pierwsze zero ma wagę **ulp**). Skoro tak szukana wartość będzie wynosić tyle, ile wybrana przez nas liczba. Dlaczego? Funkcja wyboru wartości na podstawie różnicy wygląda następująco (**Y** - liczba zaokrąglana, **M** - liczba, którą wybraliśmy):

```

M-ulp    - ulp <= Y-M <= - 1/2 ulp
M        - 1/2 ulp < Y-M < + 1/2 ulp
M+ulp    + 1/2 ulp <= Y-M <= + 1 ulp

```

Na tej podstawie można stworzyć tabelę (a raczej: odtworzyć korzystając ze wspomnianej Arytmetyki...) dla przypadków, gdy **S** = 0 (to i tak 8 przypadków i dużo pisania, a obrazują niemal całość zagadnienia):

<b>Przed zaokrągleniem</b>	...000	...001	...010	...011	...100	...101	...110	...111
<b>Po zaokrągleniu</b>	...0	...0	...0	...1	...1	...1	...1+1	...1+1

Notacja **...1+1** oznacza konieczność dodania jeszcze jedynek (nie jest to przypadek **...0**, gdyż konieczne jest jeszcze dodanie tej właśnie wartości - zwiększenie mantysy przez inkrementację). Przykład:

```

/ ----- 23 bity ----- \
Liczba 1,10101010101010101010101110
'Wybor' 1,1010101010101010101010101
Roznica: +011
Opcja: 'Wybor' + ulp
Wynik 1,1010101010101010101010101 + 1, lub inaczej:
Wynik 1,1010101010101010101010110

```

Nic dodać, nic ująć.

### Zaokrąglenie symetryczne do nieparzystej

Zaokrąglenie do nieparzystej działa podobnie do powyższego. Póki co pomijam je. Dociekliwi mogą wyprowadzić a ja się (może) wyspię.



## A. Dodatek

---

### FAQ

#### Q. Co to s<sup>1</sup> liczby pojedynczej/podwójnej precyzji?

Liczby te znajdują swój odpowiednik w C jako zmienne typu single (pojedyncza precyzja) i double (podwójna precyzja). Jak nazwa sama wskazuje - te drugie są bardziej precyzyjne, czyli w praktyce zapisywane na większej ilości bitów. Liczby pojedynczej precyzji zapisywane są na 32 bitach i mają 23 (a raczej 24, ale o tym dalej) bitową mantysę, podczas gdy liczby podwójnej precyzji mają 52 (53) bitową mantysę. By w pełni odpowiedzieć: są jeszcze liczby rozszerzone (Extended), w których mantysy mają odpowiednio 31 (32) i 63 (64) bity.

#### Q. Dlaczego 24 bity mantysy w single a nie 23?

Obie wartości są poprawne. W zapisie liczby mantysa zajmuje (przykładowo) 23 bity, jednak można (i raczej powinno się) mówić, że ma 24, gdyż tyle faktycznie bitów potrzeba do jej zapisania jeśli nie pomijamy domyślnej jedynki przed mantysą. Generalnie przy operacjach musi się też ona znaleźć w sumatorze/multiplekserze, więc powinno się mówić 24. Jednak, jak doświadczenie pokazuje, co ciało pedagogiczne to zwyczaj, a te ciała wola 23.

#### Q. Po co liczby zdenormalizowane?

Trudno powiedzieć. Być może, by uwzględnić szczególne przypadki wymagające większej precyzji? A może dlatego, że zostały wolne kombinacje Mantysa/Wykładnik? Jaki jest powód nie ma znaczenia. Liczby te są w standardzie i stały się przyczyną opóźnienia w jego zatwierdzeniu. Do dziś wiele układów arytmetyki zmiennoprzecinkowej nie ma bezpośredniej implementacji tych liczb.

#### Q. Czy są inne standardy zmiennoprzecinkowe?

Tak i to niestety kilka. Jest pokrecony zapis logarytmiczny i oczywiście nowsza wersja standardu IEEE - IEEE 854. Ten ostatni implikuje między innymi konieczność implementacji liczb zmiennoprzecinkowych o bazie 10 (nie tylko 2).

#### Q. Jak, czy i kiedy dokonujemy wyboru metody zaokrąglania?

Niestety wszystko zależy od jednostki, która obliczenia wykonuje. Standard nie określa gdzie mają znajdować się informacje o aktualnie wykorzystywanym typie z. Jedni producenci umieszczają bity charakterystyczne w słowie statusowym mikrokontrolera, inni w specjalnych rejestrach, jeszcze inni w ogóle nie dają możliwości wyboru.

### Kontakt z autorami

email: [IronHand](#)

email: [Ryan O'Connell](#)