

**ОСОБЕННОСТИ АРИФМЕТИКИ  
С ПЛАВАЮЩЕЙ ЗАПЯТОЙ  
В СОВРЕМЕННЫХ КОМПЬЮТЕРАХ**

. 60- 70-

. -

. -

, -

. -

, -

. -

, -

. -

, -

. -

, -

. -

70- 80-

, -

IEEE 754. -

, 80- -

, -

, -



( $|M| = 3$ ).  
 $101_2$ .

$2^1 = 2,$

$2^0 = 1.$   
 $E = 2.$

$2^2 = 4,$

«1.01e + 2».

$E = 1.$

$1.01e + 1 = 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1} = 2 + 0,5 = 2,5.$

$|M| = 4.$  «2»

$2 = 10 ( ) = 1.000e + 1 = 0.100e + 2 = 0.010e + 3.$

( . 2),

IEEE 754,

– IEEE 754-2008 [6].



. 2.

IEEE-754

( $s$ ), ( $M$ ) ( $E$ ) :

$(-1)^s \times 1.M \times 2^E.$

(single).

(half), (double) (extended)

(float/single)

8

– 23.

E-127.

IEEE 754 «0»  
 $E = E_{min} - 1$  ( single – 127)

IEEE 754

(±) (NaN).  
 $E = E_{max} + 1$

(, 3/0 = +, - 3/0 = -),  
 0/0

NaN («not a number»)

IEEE 754 NaN  
 $E = E_{max} + 1$ , NaN NaN.

NaN NaN,

NaN

“+(-)”, “0 ×”, “0/0, /”, “√x” (x < 0).  
 «0». IEEE754 «0»

«0», 3•(+ 0) = + 0,  
 3•(- 0) = - 0.

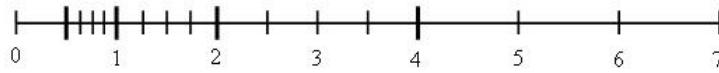
1/ 1/ - 0.  $1/(1/x) = x$ ,  $x = \pm$ ,  
 «+ 0 = - 0».

: «(+ /0) + = + », «(+ / - 0) + = NaN».  
 NaN, NaN, NaN.

«1/ = 0». x = 1, 2, 3, 4:

$$F(x) = 7 - \frac{3}{x-2 - \frac{1}{x-7 + \frac{10}{x-2 - \frac{2}{x-3}}}}$$

$$|M| = 2^{-1} E - 2 \quad (3)$$

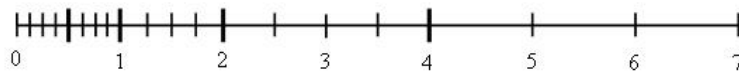


3.

1,00. (0 - 0,5), 0,5 1  
 (0,5 - 0,625), 0,5 0  
 0, «1,5 - 1,25 = 0» (3).  
 IEEE 754  
 $E = E_{min} - 1$  (float «-127»)  
 $E = E_{min}$ ,

$$(-1)^s \times 1.M \times 2^E, \quad E_{min} \leq E \leq E_{max} \quad (4)$$

$$(-1)^s \times 0.M \times 2^{E_{min}}, \quad E = E_{min} - 1 \quad (E_{min} = -1, E = -2,$$

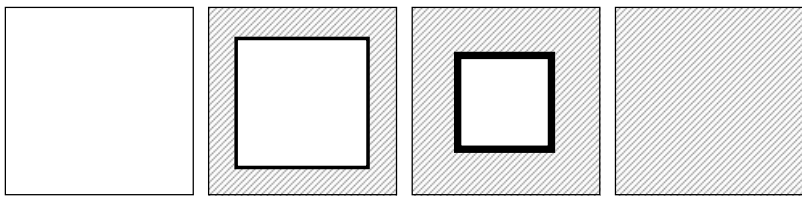


4.

0 0,5 (0,5 - 0,25)  
 1,5 - 1,25). [7],

[7]

,  
:  
,  
,  
-  
.5  
,  
( .5, ).  
,  
.5,  
-  
( .5, ),  
,  
( .5, ).

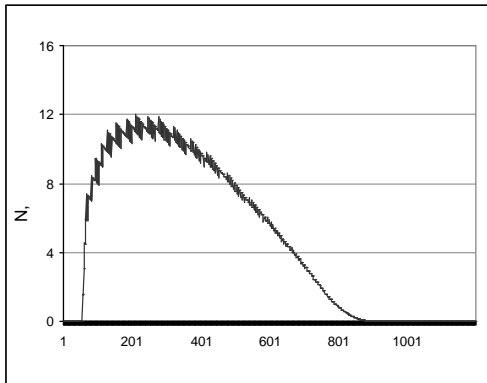
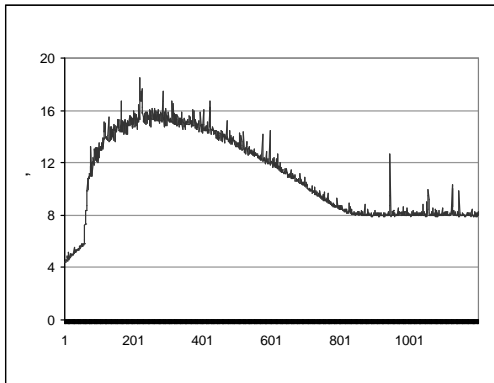


.5. : - -  
; -  
; -

500 x 500, Intel Core 2 Duo 2.14  
V = 100,

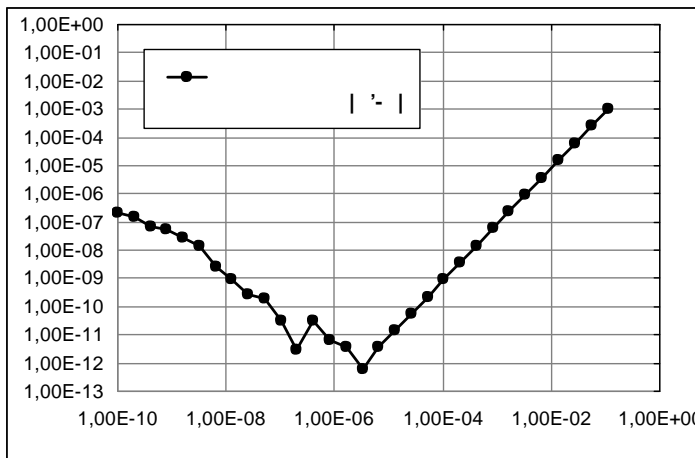
.6, ,  
~ 6 %

30,  
,  
,  
( 1000- )



. 6. : - ; -

(1). . 7



$$\int_0^1 \frac{4}{1-x^2} = \pi(1)$$

. 7.

IEEE754

{n}

$$n < n + 1 \Rightarrow F(n) < F(n + 1),$$

F(n) -

n,

```

float a=0.5;
int n = *((int*) &a);
float b = *((float*) &(&n));
printf("      %e      : %e,      (%e)\n", a, b, b-a);

```

[8].

: « $(x - y)(x + y)$ » [5, 8]

$$(1020 + 1) - 1020 = 0 \quad (1020 - 1020) + 1 = 1.$$

```

double s = 0.0;
for (int i=0; i<n; i++) s = s + t[i];

```

, n 2): ( ) ( -



```

double sa[2], s;
sa[0]=sa[1]=0.0;
for (int i=0; i<n/2; i++) {
    sa[0]=sa[0]+t[i*2+0];
    sa[1]=sa[1]+t[i*2+1];
}
s=sa[0]+sa[1];

```

«?»

$x$	$y$	$x < y ? x: y$	$x \leq y ? x: y$	$x > y ? y: x$	$x \geq y ? y: x$
+0	-0	-0	+0	+0	-0
NaN	1	1	1	NaN	NaN

±0 NaN,

float-

```

float fValue = 0.2;
if (fValue == 0.2) DoStuff();

```

fValue - 0,2 -

```

if (fabs(fValue - fExpected) < 0.0001) DoStuff();

```

«10000»,

(10000,000977).

[9].

---

```

bool AlmostEqual (float A, float B, int maxUlp)
{
    // maxUlp
    // NaN
    assert(maxUlp > 0 && maxUlp < 4 * 1024 * 1024);
    int aInt = *(int*)&A;
    // aInt,
    //
    //
    if (aInt < 0) aInt = 0x80000000 - aInt;
    // bInt
    int bInt = *(int*)&B;
    if (bInt < 0) bInt = 0x80000000 - bInt;
    unsigned int intDiff = abs(aInt - bInt);
    if (intDiff <= maxUlp)
        return true;
    return false;
}

```

$\text{maxUlp}$  («Units-In-Last-Place») –  
 ,  
 ( , ),  
 ,  $\text{maxUlp} = 16$ , , 4 ( $\log_2 16$ )  
 , 10000 , 0,0146,  
 0.001, 0.00000001 (10 – 8).

( , )  
 , [10],  
 20 ,  
 ,  

$$S = \sum_{i=1}^n x_i$$

$$S = S_1 + S_2 = \sum_{i=1}^{n/2} x_i + \sum_{i=n/2+1}^n x_i.$$

```

[8]
(
),
» [11],

//
double CompensatedSum(double *x, int count)
{
    s=0; e=0; temp=0; y=0;
    for (i=0; i<count; i++) {
        temp = s;
        y = x[i] + e;
        s = temp + y;
        e = (temp - s) + y;
    }
}

```

PLASMA .

[12].

[13].

IEEE 754,

---

*R. Iushchenko, O. Iushchenko*

#### SPECIFICITIES OF FLOATING POINT ARITHMETIC ON MODERN COMPUTERS

An overview of modern floating point arithmetic is presented. The aspects, which are different from real numbers arithmetic, and, thus, cause a lot of numeric problems are outlined. These problems become substantial on large amounts of data and in parallel computations. Nevertheless, it is possible to guarantee the required accuracy of the results without considerable loss of performance.

1. . . . . // . - 2004. - 6. - . 65 - 72.
2. . . . . // . - 2009. - 6. - . 88 - 96.
3. <http://www.cs.unc.edu/~ibr/projects/paranoia/>.
4. <http://software.intel.com/en-us/videos/tim-mattson-floating-points-arent-real/>.
5. *Goldberg D.* What Every Computer Scientist Should Know About Floating-Point Arithmetic // ACM Computing Surveys. - 1991. - Vol. 1, N 23. - P. 5 - 48.
6. *IEEE 754-2008 Standard for Floating-Point Arithmetic.*
7. *Bjørndalen J.M., Anshus O.J.* Trusting floating point benchmarks - are your benchmarks really data independent?, Proceedings of the 8th international conference on Applied parallel computing: state of the art in scientific computing, June 18 - 21, 2006, Umeå, Sweden.
8. *Higham N.* Accuracy and Stability of Numerical Algorithms. - Philadelphia: 2002, SIAM. - . 680.

- 
9. *Bruce Dawson*,  
<http://www.cygnus-software.com/papers/comparingfloats/comparingfloats.htm>
  10. *Tim Mattson*,  
<http://software.intel.com/en-us/videos/tim-mattson-floating-points-arent-real/>.
  11. *Kahan W.* Implementation of algorithms. Technical Report 20, Department of Computer Science, University of California, Berkeley, CA, USA, 1973.
  12. *Buttari A., Dongarra J., Kurzak J., Luszczek P., Tomov S.* "Using Mixed Precision for Sparse Matrix Computations to Enhance the Performance while Achieving 64-bit Accuracy," ACM Transactions on Mathematical Software. – 2008. – Vol. 34, N 4. – . 17 – 22.
  13. . . . .  
 // . – 2009. – 6. –  
 . 172 – 176.

15.03.2016

**Об авторах:**