

Продолжаем говорить про арифметику

Андрей Акиньшин, JetBrains

DotNext 2016 Moscow

Знать математику — полезно



О чём будем разговаривать?

Сегодня в выпуске:

- Весёлые истории про арифметические баги
- Хитрые битовые хаки
- Удивительные рантаймы и JIT-компиляторы
- Коварные баги в .NET
- Неожиданные ASM-листинги
- Загадочный IEEE 754 и денормализованные числа

О чём будем разговаривать?

Сегодня в выпуске:

- Весёлые истории про арифметические баги
- Хитрые битовые хаки
- Удивительные рантаймы и JIT-компиляторы
- Коварные баги в .NET
- Неожиданные ASM-листинги
- Загадочный IEEE 754 и денормализованные числа
- **Проблемы с точностью:
кто виноват и что делать?**
- **Скорость вычислений:
от чего зависит и как улучшать?**

Отказ от ответственности

Нужно понимать:

- Все бенчмарки вымыщены, любые совпадения случайны
- Приведённые примеры нужны только для иллюстрации идей
- Погрешности малы, для целей доклада ими можно пренебречь
- На вашем железе цифры могут быть другие, это нормально

Конфигурация для всех примеров:

- BenchmarkDotNet v0.10.1 (supported by the .NET Foundation)
- Haswell Core i7-4702MQ CPU 2.20GHz, Windows 10
- .NET Framework 4.6.2 + clrjit/compatjit-v4.6.1586.0
- Mono 4.6.2 (x86/x64)
- Gist с исходниками: <https://git.io/v1lUV>

Quake 3 and Fast Inverse Square Root

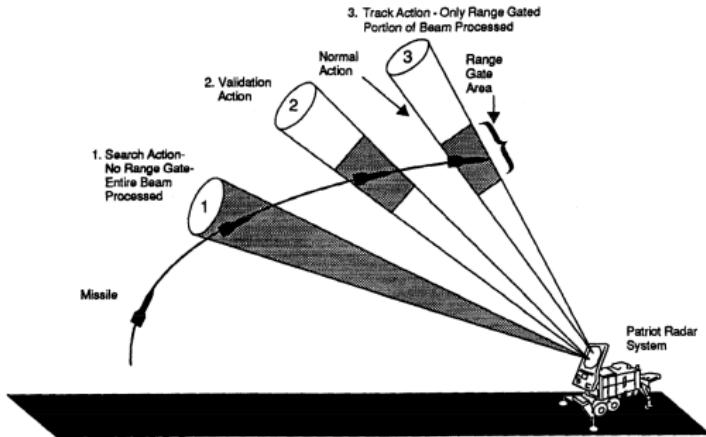


Quake 3 and Fast Inverse Square Root



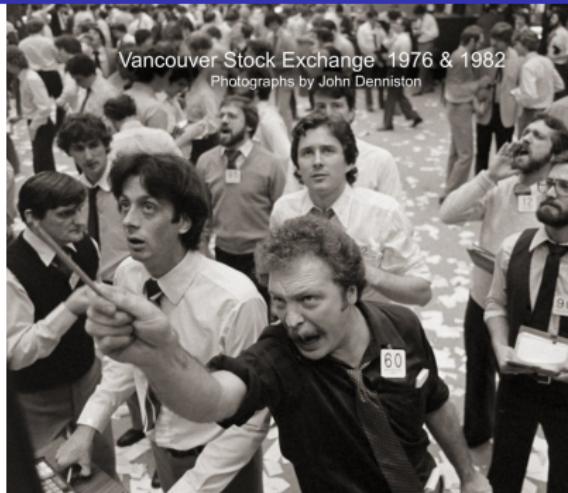
```
float Q_rsqrt(float number) { // 1/sqrt(number)
    long i; float x2, y;
    const float threehalfs = 1.5F;
    x2 = number * 0.5F; y = number;
    i = *(long*)&y; // evil floating point bit level hacking
    i = 0x5f3759df - (i >> 1); // what the fuck?
    y = *(float*)&i;
    y = y * (threehalfs - (x2 * y * y)); // 1st iteration
    return y;
}
```

The Patriot Missile Failure

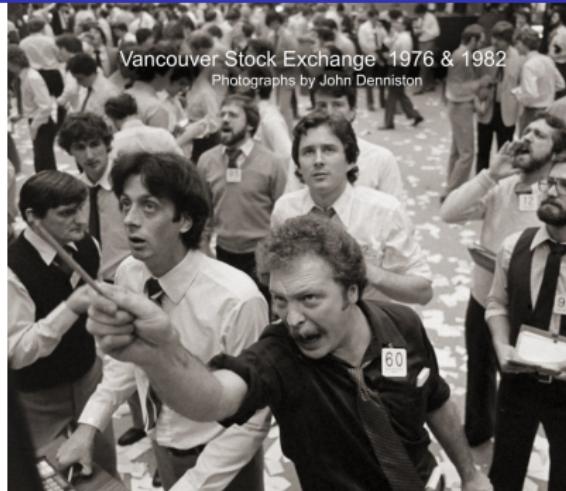


- Война в Персидском заливе, 25 февраля 1991
- Ракета Patriot не смогла сбить иракский снаряд Scud из-за бага
- Скорость Scud: ≈ 3750 миль в час
- Время хранилось в 24-битном регистре с точностью до 0.1 сек
- Ошибка на 1 тик составляет $\approx 9.5 \cdot 10^{-8}$ сек
- Суммарная ошибка за 100 часов: 0.3433 сек
- Итоговая погрешность > 0.5 км
- 28 человек убиты, 98 ранены, см. также GAO/IMTEC-92-26

Vancouver Stock Exchange

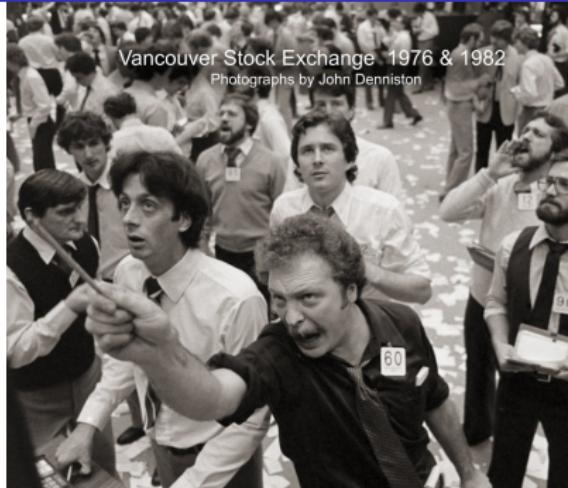


Vancouver Stock Exchange



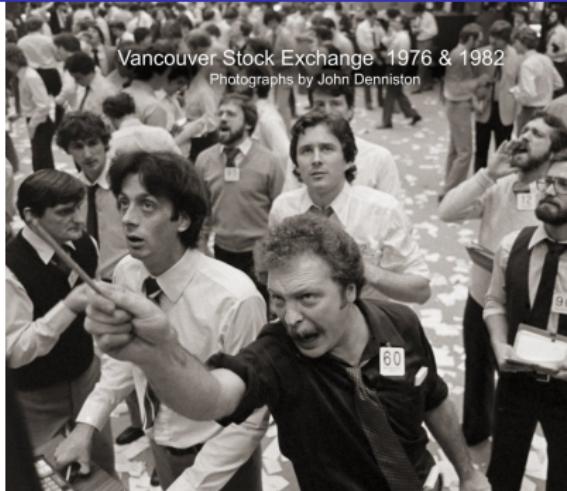
- Январь 1982: введён индекс = 1000

Vancouver Stock Exchange



- Январь 1982: введён индекс = 1000
- После каждой транзакции индекс обновлялся и округлялся до 3 знаков после точки (≈ 2800 раз в день)

Vancouver Stock Exchange

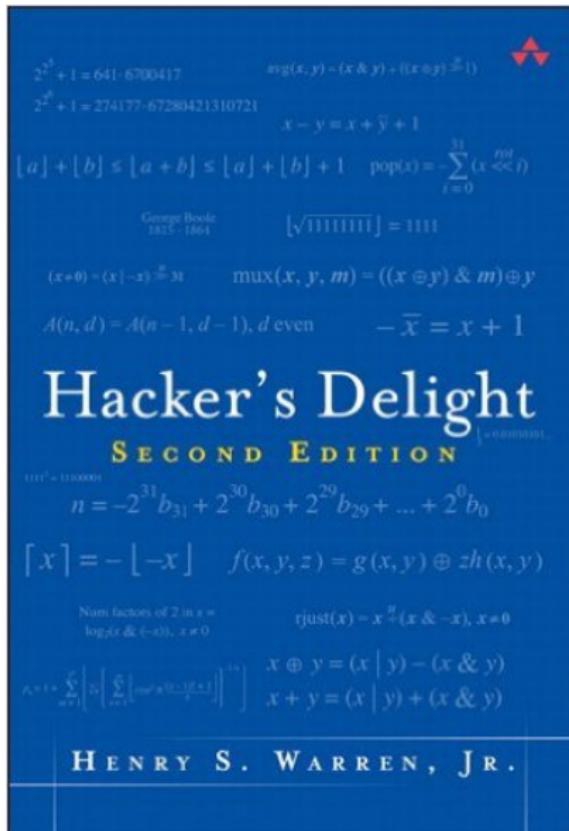


Vancouver Stock Exchange 1976 & 1982

Photographs by John Denniston

- Январь 1982: введён индекс = 1000
- После каждой транзакции индекс обновлялся и округлялся до 3 знаков после точки (≈ 2800 раз в день)
- Ноябрь 1983: ребята скорректировали ошибки
Изменение индекса: **524.811 → 1098.892**

Поговорим про битовые хаки



Ищем минимум

Задача. Дано два числа a и b типа `int`.

Необходимо реализовать функцию `Min`, которая вернёт наименьшее из этих двух чисел.

Задача. Дано два числа a и b типа `int`.

Необходимо реализовать функцию `Min`, которая вернёт наименьшее из этих двух чисел.

Решение А.

```
int MinA(int x, int y) => Math.Min(x, y);
```

Задача. Дано два числа a и b типа `int`.

Необходимо реализовать функцию `Min`, которая вернёт наименьшее из этих двух чисел.

Решение А.

```
int MinA(int x, int y) => Math.Min(x, y);
```

Решение В.

```
int MinB(int x, int y) => x < y ? x : y;
```

Задача. Дано два числа a и b типа `int`.

Необходимо реализовать функцию `Min`, которая вернёт наименьшее из этих двух чисел.

Решение А.

```
int MinA(int x, int y) => Math.Min(x, y);
```

Решение В.

```
int MinB(int x, int y) => x < y ? x : y;
```

Решение С.

```
int MinC(int x, int y) =>
    x & ((x - y) >> 31) |
    y & (~ (x - y) >> 31);
```

(Объяснение) Ищем минимум

```
int MinC(int x, int y) =>  
    x & ( (x - y) >> 31) |  
    y & (~ (x - y) >> 31);
```

x = 8	0	0	...	1	0	0	0
y = 3	0	0	...	0	0	1	1
x-y	0	0	...	0	1	0	1
(x-y)>>31	0	0	...	0	0	0	0
~(x-y)>>31	1	1	...	1	1	1	1
x&((x-y)>>31)	0	0	...	0	0	0	0
y&(~(x-y)>>31)	0	0	...	0	0	1	1
Result	0	0	...	0	0	1	1

(Бенчмарки) Ищем минимум

```
const int N = 100001; int[] a, b, c; // int[N] arrays
[Benchmark] void MinB() {
    for (int i = 0; i < N; i++) {
        int x = a[i], y = b[i];
        c[i] = x < y ? x : y;
    }
}
[Benchmark] void MinC() {
    for (int i = 0; i < N; i++) {
        int x = a[i], y = b[i];
        c[i] = x & ((x - y) >> 31) | y & (~ (x - y) >> 31);
    }
}
```

(Бенчмарки) Ищем минимум

```
const int N = 100001; int[] a, b, c; // int[N] arrays
[Benchmark] void MinB() {
    for (int i = 0; i < N; i++) {
        int x = a[i], y = b[i];
        c[i] = x < y ? x : y;
    }
}
[Benchmark] void MinC() {
    for (int i = 0; i < N; i++) {
        int x = a[i], y = b[i];
        c[i] = x & ((x - y) >> 31) | y & (~ (x - y) >> 31);
    }
}
```

	Random	Const	MinB	MinC
	MinB	MinC	MinB	MinC
LegacyJIT-x86	≈782μs	≈292μs	≈160μs	≈224μs
LegacyJIT-x64	≈525μs	≈157μs	≈71μs	≈124μs
RyuJIT-x64	≈602μs	≈248μs	≈221μs	≈292μs
Mono4.6.2-x64	≈255μs	≈349μs	≈212μs	≈290μs

(ASM) Ищем минимум

```
; RyuJIT-x641
cmp        ecx, edx
jl         LESS
mov        eax, edx
ret
LESS:
mov        eax, ecx
ret
```

¹Roslyn 1.3.1.60616, /o+

(ASM) Ищем минимум

```
; RyuJIT-x641
cmp        ecx, edx
jl         LESS
mov        eax, edx
ret
LESS:
mov        eax, ecx
ret
```

```
; Mono4.6.2-x64
sub       $0x18,%rsp
mov        %rsi,(%rsp)
mov        %rdi,0x8(%rsp)
mov        %rcx,%rdi
mov        %rdx,%rsi
cmp        %esi,%edi
mov        %rsi,%rax
cmovl    %rdi,%rax
mov        (%rsp),%rsi
mov        0x8(%rsp),%rdi
add       $0x18,%rsp
retq
```

cmovl (0F 4C): Conditional move (if less)

¹Roslyn 1.3.1.60616, /o+

Делим на константу

Задача. Дано число `n` типа `int`. Необходимо реализовать функцию `Div3`, которая делит заданное число на 3.

Делим на константу

Задача. Дано число n типа `int`. Необходимо реализовать функцию `Div3`, которая делит заданное число на 3.

Решение А.

```
int Div3A(int n) =>
    n / 3;
```

Делим на константу

Задача. Дано число n типа `int`. Необходимо реализовать функцию `Div3`, которая делит заданное число на 3.

Решение А.

```
int Div3A(int n) =>
    n / 3;
```

Решение В.

```
int Div3B(int n) =>
    (int)((n * 0xAAAAAAAB) >> 33);
```

Делим на константу

Задача. Дано число n типа `int`. Необходимо реализовать функцию `Div3`, которая делит заданное число на 3.

Решение А.

```
int Div3A(int n) =>
    n / 3;
```

Решение В.

```
int Div3B(int n) =>
    (int)((n * 0xAAAAAAAB) >> 33);
```

$$0xAAAAAAAB_{16} = \frac{2^{33}}{3} + 1$$

(Бенчмарки) Делим на константу

```
int N = 100001; int[] a, b; // random arrays
[Benchmark] void Div3A() {
    for (int i = 0; i < N; i++)
        b[i] = a[i] / 3;
}
[Benchmark] void Div3B() {
    for (int i = 0; i < N; i++)
        b[i] = (int)((a[i] * 0xAAAAAAAB) >> 33);
}
```

(Бенчмарки) Делим на константу

```
int N = 100001; int[] a, b; // random arrays
[Benchmark] void Div3A() {
    for (int i = 0; i < N; i++)
        b[i] = a[i] / 3;
}
[Benchmark] void Div3B() {
    for (int i = 0; i < N; i++)
        b[i] = (int)((a[i] * 0xAAAAAAAB) >> 33);
}
```

	Div3A	Div3B
LegacyJIT-x86	≈376μs	≈334μs
LegacyJIT-x64	≈122μs	≈67μs
RyuJIT-x64	≈173μs	≈140μs
Mono4.6.2-x86	≈366μs	≈1 139μs
Mono4.6.2-x64	≈446μs	≈176μs

(Объяснение) Делим на константу

```
long Mul(int x, uint y) => x * y;
```

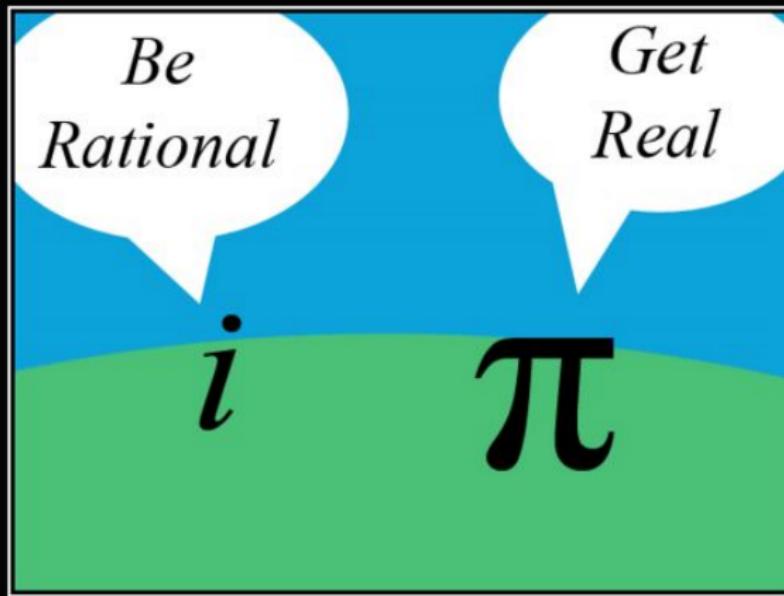
(Объяснение) Делим на константу

```
long Mul(int x, uint y) => x * y;
```

```
; Mono-x86
gram_Mul:
push %rbp
mov %esp,%ebp
push %rdi
sub $0x34,%esp
cmpl $0xffffffff,0x8(%rbp)
setg %al
movzbl %al,%eax
mov 0x8(%rbp),%rcx
mov %ecx,-0x10(%rbp)
mov %eax,-0xc(%rbp)
mov 0xc(%rbp),%eax
mov %eax,-0x18(%rbp)
movl $0x0,-0x14(%rbp)
mov -0x14(%rbp),%eax
mov %eax,0xc(%rsp)
mov -0x18(%rbp),%eax
mov %eax,0x8(%rsp)
mov -0xc(%rbp),%eax
mov %eax,0x4(%rsp)
```

```
mov -0x10(%rbp),%eax
mov %eax,(%rsp)
callq c2d080c gram_Mul+0xc2d080c
mov %edx,-0xc(%rbp)
mov %eax,-0x10(%rbp)
mov 0x102ebc54(%rip),%eax
mov -0x10(%rbp),%ecx
mov %ecx,-0x18(%rbp)
mov -0xc(%rbp),%ecx
mov %ecx,-0x14(%rbp)
test %eax,%eax
jne gram_Mul+0xa8
jmp gram_Mul+0x72
mov -0x20(%rbp),%eax
mov %eax,-0x18(%rbp)
mov -0x1c(%rbp),%eax
mov %eax,-0x14(%rbp)
mov -0x18(%rbp),%eax
mov %eax,-0x10(%rbp)
mov -0x14(%rbp),%eax
; И ёщё куча строк...
```

Нецелые числа — это сложно



MATH JOKES

If you get them, you probably don't have any friends.

Обычная арифметика не работает

$$(a + b) + c \neq a + (b + c)$$

$$(a \cdot b) \cdot c \neq a \cdot (b \cdot c)$$

$$(a + b) \cdot c \neq a \cdot c + b \cdot c$$

$$a^{x+y} \neq a^x \cdot a^y$$

...

Обычная арифметика не работает

Жить сложно:

$$0.1d + 0.2d \neq 0.3d$$

$$0.1d \approx 0.10000000000000005551115123125783$$

$$+ 0.2d \approx 0.200000000000000011102230246251565$$

$$0.300000000000000044408920985006262$$

$$0.3d \approx 0.29999999999999988897769753748435$$

Обычная арифметика не работает

Жить сложно:

$$0.1d + 0.2d \neq 0.3d$$

$$\begin{array}{r} 0.1d \approx 0.10000000000000005551115123125783 \\ + 0.2d \approx 0.200000000000000011102230246251565 \\ \hline 0.300000000000000044408920985006262 \\ 0.3d \approx 0.29999999999999988897769753748435 \end{array}$$

Хорошая практика:

```
bool AreEqual(double a, double b,  
              double eps = 1e-9)  
{  
    return Math.Abs(a - b) < eps;  
}
```

Поговорим про конвертацию

Задача. Что выведет следующий код?

```
double d1 = 0.84551240822557006;  
string str = d1.ToString("R");1  
double d2 = double.Parse(str);  
WriteLine(d1 == d2);
```

Поговорим про конвертацию

Задача. Что выведет следующий код?

```
double d1 = 0.84551240822557006;  
string str = d1.ToString("R");1  
double d2 = double.Parse(str);  
WriteLine(d1 == d2);
```

MSDN: Standard Numeric Format Strings

¹ The round-trip ("R") format specifier is used to ensure that a numeric value that is converted to a string will be parsed back into the same numeric value.

Поговорим про конвертацию

Задача. Что выведет следующий код?

```
double d1 = 0.84551240822557006;
string str = d1.ToString("R");1
double d2 = double.Parse(str);
WriteLine(d1 == d2);
```

MSDN: Standard Numeric Format Strings

¹ The round-trip ("R") format specifier is used to ensure that a numeric value that is converted to a string will be parsed back into the same numeric value.

.NET-x64	.NET-x86	Mono
false (баг)	true	true

Загадка на производительность

Вопрос. Какая программа работает намного медленнее остальных под LegacyJIT-x86?

```
const int N = 100001;
public double Sum()
{
    double x = 1, y = 1;

    for (int i = 0; i < N; i++)
        x = x + y;

    return x;
}
// A
Sum();
// B
Write(Stopwatch.Frequency);
Sum();
```

Загадка на производительность

Вопрос. Какая программа работает намного медленнее остальных под LegacyJIT-x86?

```
const int N = 100001;
public double Sum()
{
    double x = 1, y = 1;

    for (int i = 0; i < N; i++)
        x = x + y;

    return x;
}
// A
Sum();
// B
Write(Stopwatch.Frequency);
Sum();
```

```
const int N = 100001;
public double Sum2()
{
    double x = 1, y = 1;
    var sw = Stopwatch.StartNew();
    for (int i = 0; i < N; i++)
        x = x + y;
    sw.Stop();
    return x;
}
// C
Sum2();
// D
Write(Stopwatch.Frequency);
Sum2();
```

Загадка на производительность

Вопрос. Какая программа работает намного медленнее остальных под LegacyJIT-x86?

```
const int N = 100001;
public double Sum()
{
    double x = 1, y = 1;

    for (int i = 0; i < N; i++)
        x = x + y;

    return x;
}
// A
Sum();
// B
Write(Stopwatch.Frequency);
Sum();
```

```
const int N = 100001;
public double Sum2()
{
    double x = 1, y = 1;
    var sw = Stopwatch.StartNew();
    for (int i = 0; i < N; i++)
        x = x + y;
    sw.Stop();
    return x;
}
// C
Sum2();
// D
Write(Stopwatch.Frequency);
Sum();
```

Ответ. Программа С:

	A	B	C	D
LegacyJIT-x86	$\approx 100\mu s$	$\approx 100\mu s$	$\approx 335\mu s$	$\approx 100\mu s$

Смотрим ASM

```
// D
// Статический конструктор класса
// Stopwatch вызывается при
// обращении к Stopwatch.Frequency
Write(Stopwatch.Frequency);

// Поэтому нам не нужно
// вызывать его из метода.

Sum2();
```

Смотрим ASM

```
// D
// Статический конструктор класса
// Stopwatch вызывается при
// обращении к Stopwatch.Frequency
Write(Stopwatch.Frequency);

// Поэтому нам не нужно
// вызывать его из метода.
Sum2();
```

```
// C
// Статический конструктор класса
// Stopwatch ещё не вызывался.

// Поэтому нам придётся
// вызвать его из метода.
Sum2();
```

Смотрим ASM

```
// D  
// Статический конструктор класса  
// Stopwatch вызывается при  
// обращении к Stopwatch.Frequency  
Write(Stopwatch.Frequency);  
  
// Поэтому нам не нужно  
// вызывать его из метода.  
Sum2();
```

```
// C  
// Статический конструктор класса  
// Stopwatch ещё не вызывался.  
  
// Поэтому нам придётся  
// вызвать его из метода.  
Sum2();
```

```
; x + y (A, B, D)  
fld1  
faddp st(1),st
```

```
; x + y (C)  
fld1  
fadd qword ptr [ebp-0Ch]  
fstp qword ptr [ebp-0Ch]
```

Browse Standards > IEEE Std 754-2008 ... 

754-2008 - IEEE Standard for Floating-Point Arithmetic



Full Text

Sign-In or Purchase

Revision of ANSI/IEEE Std 754-1985

Status: Active - Approved

Abstract

References

Versions

Definitions

Cited By

Keywords



Download
Citations



Email



Print



Request
Permissions



Export



This standard specifies interchange and arithmetic formats and methods for binary and decimal floating-point arithmetic in computer programming environments. This standard specifies exception conditions and their default handling. An implementation of a floating-point system conforming to this standard may be realized entirely in software, entirely in hardware, or in any combination of software and hardware. For operations specified in the normative part of this standard, numerical results and exceptions are uniquely determined by the values of the input data, sequence of operations, and destination formats, all under user control.

Date of Publication :

Aug. 29 2008

INSPEC Accession Number:

10369888

Status :

Active

DOI:

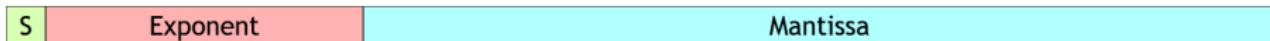
10.1109/IEEESTD.2008.4610935

Нормализованные числа



$$V = (-1)^S \cdot 1.M \cdot 2^{E-E_{bias}}, \quad E_{min} \leq E \leq E_{max}$$

Нормализованные числа



$$V = (-1)^S \cdot 1.M \cdot 2^{E-E_{bias}}, \quad E_{min} \leq E \leq E_{max}$$

	Sign	Exponent	Mantissa	E_{bias}
float	1	8	23	127
double	1	11	52	1023
80bit	1	15	64	16383

Нормализованные числа

S	Exponent	Mantissa
---	----------	----------

$$V = (-1)^S \cdot 1.M \cdot 2^{E-E_{bias}}, \quad E_{min} \leq E \leq E_{max}$$

	Sign	Exponent	Mantissa	E_{bias}
float	1	8	23	127
double	1	11	52	1023
80bit	1	15	64	16383

	Digits	Lower	Upper
float	≈ 7.2	$1.2 \cdot 10^{-38}$	$3.4 \cdot 10^{+38}$
double	≈ 15.9	$2.3 \cdot 10^{-308}$	$1.7 \cdot 10^{+308}$
80bit	≈ 19.2	$3.4 \cdot 10^{-4932}$	$1.1 \cdot 10^{+4932}$

Пример

0	1	0	0	1	1	1	0	0	1	1	0	1	1	0	0	1	1	0	1	0	1	0	0	1	0	1	0	0	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$$S = 0$$

$$E = 10011100_2 = 156_{10}$$

$$M = 1.11011100110101100101001_2 = 15\,625\,001 \cdot 2^{-23}$$

$$\begin{aligned}V &= (-1)^0 \cdot \left(15\,625\,001 \cdot 2^{-23}\right) \cdot 2^{156-127} \\&= 15\,625\,001 \cdot 2^6 = 1\,000\,000\,064\end{aligned}$$

Пример

0	1	0	0	1	1	1	0	0	1	1	1	0	0	1	1	0	1	0	1	1	0	0	1	0	1	0	0	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$$S = 0$$

$$E = 10011100_2 = 156_{10}$$

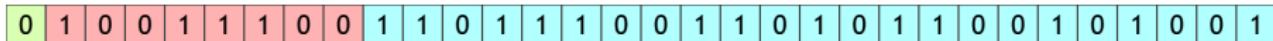
$$M = 1.11011100110101100101001_2 = 15\,625\,001 \cdot 2^{-23}$$

$$\begin{aligned}V &= (-1)^0 \cdot \left(15\,625\,001 \cdot 2^{-23}\right) \cdot 2^{156-127} \\&= 15\,625\,001 \cdot 2^6 = 1\,000\,000\,064\end{aligned}$$

Задача. Чему равно значение выражения?

```
((float)1000000064).ToString("G10")
```

Пример



$$S = 0$$

$$E = 10011100_2 = 156_{10}$$

$$M = 1.11011100110101100101001_2 = 15\,625\,001 \cdot 2^{-23}$$

$$\begin{aligned}V &= (-1)^0 \cdot \left(15\,625\,001 \cdot 2^{-23}\right) \cdot 2^{156-127} \\&= 15\,625\,001 \cdot 2^6 = 1\,000\,000\,064\end{aligned}$$

Задача. Чему равно значение выражения?

```
((float)1000000064).ToString("G10")
```

Решение.

```
1000000060 // 9 значащих цифр
```

Денормализованные числа



$$V = (-1)^S \cdot 0.M \cdot 2^{E_{denorm}}, \quad E_{denorm} = E_{bias} - 1$$

Денормализованные числа



$$V = (-1)^S \cdot 0.M \cdot 2^{E_{denorm}}, \quad E_{denorm} = E_{bias} - 1$$

	E_{denorm}	Lower	Upper
float	-126	$1.4 \cdot 10^{-45}$	$1.2 \cdot 10^{-38}$
double	-1022	$5.0 \cdot 10^{-324}$	$2.3 \cdot 10^{-308}$
80bit	-16382	$1.9 \cdot 10^{-4951}$	$3.4 \cdot 10^{-4932}$

Денормализованные числа



$$V = (-1)^S \cdot 0.M \cdot 2^{E_{denorm}}, \quad E_{denorm} = E_{bias} - 1$$

	E_{denorm}	Lower	Upper
float	-126	$1.4 \cdot 10^{-45}$	$1.2 \cdot 10^{-38}$
double	-1022	$5.0 \cdot 10^{-324}$	$2.3 \cdot 10^{-308}$
80bit	-16382	$1.9 \cdot 10^{-4951}$	$3.4 \cdot 10^{-4932}$

- Хорошие новости: $a \neq b \Rightarrow a - b \neq 0$

Денормализованные числа



$$V = (-1)^S \cdot 0.M \cdot 2^{E_{denorm}}, \quad E_{denorm} = E_{bias} - 1$$

	E_{denorm}	Lower	Upper
float	-126	$1.4 \cdot 10^{-45}$	$1.2 \cdot 10^{-38}$
double	-1022	$5.0 \cdot 10^{-324}$	$2.3 \cdot 10^{-308}$
80bit	-16382	$1.9 \cdot 10^{-4951}$	$3.4 \cdot 10^{-4932}$

- Хорошие новости: $a \neq b \Rightarrow a - b \neq 0$
- Плохие новости: денормализованные числа работают медленно ☹

(Задачка) Денормализованные числа

Задача. Подсчитать $0.96^{100\,000}$, $0.96^{1\,000\,000}$, используя **double**.

(Задачка) Денормализованные числа

Задача. Подсчитать $0.96^{100\,000}$, $0.96^{1\,000\,000}$, используя **double**.

Решение А.

```
public double PowerA() {  
    double res = 1.0;  
    for (int i = 0; i < N; i++)  
        res = res * 0.96;  
    return res;  
}
```

(Задачка) Денормализованные числа

Задача. Подсчитать $0.96^{100\,000}$, $0.96^{1\,000\,000}$, используя **double**.

Решение А.

```
public double PowerA() {  
    double res = 1.0;  
    for (int i = 0; i < N; i++)  
        res = res * 0.96;  
    return res;  
}
```

Решение В.

```
private double resB;  
public double PowerB() {  
    resB = 1.0;  
    for (int i = 0; i < N; i++)  
        resB = resB * 0.96;  
    return resB;  
}
```

(Задачка) Денормализованные числа

Задача. Подсчитать $0.96^{100\,000}$, $0.96^{1\,000\,000}$, используя `double`.

Решение А.

```
public double PowerA() {  
    double res = 1.0;  
    for (int i = 0; i < N; i++)  
        res = res * 0.96;  
    return res;  
}
```

Решение В.

```
private double resB;  
public double PowerB() {  
    resB = 1.0;  
    for (int i = 0; i < N; i++)  
        resB = resB * 0.96;  
    return resB;  
}
```

Решение С.

```
public double PowerC() {  
    double res = 1.0;  
    for (int i = 0; i < N; i++)  
        res = res * 0.96  
            + 0.1 - 0.1;  
    return res;  
}
```

```
[Params(100000, 1000000)]  
public int N;
```

(Бенчмарки) Денормализованные числа

	N	PowerA	PowerB	PowerC
LegacyJIT-x86	10^5	$\approx 168\mu s$	$\approx 19\ 264\mu s$	$\approx 370\mu s$
RyuJIT-x64	10^5	$\approx 3\ 961\mu s$	$\approx 4\ 142\mu s$	$\approx 370\mu s$
LegacyJIT-x86	10^6	$\approx 152\ 770\mu s$	$\approx 227\ 781\mu s$	$\approx 3\ 851\mu s$
RyuJIT-x64	10^6	$\approx 47\ 956\mu s$	$\approx 49\ 670\mu s$	$\approx 3\ 782\mu s$

(Бенчмарки) Денормализованные числа

	N	PowerA	PowerB	PowerC
LegacyJIT-x86	10^5	$\approx 168\mu s$	$\approx 19\,264\mu s$	$\approx 370\mu s$
RyuJIT-x64	10^5	$\approx 3\,961\mu s$	$\approx 4\,142\mu s$	$\approx 370\mu s$
LegacyJIT-x86	10^6	$\approx 152\,770\mu s$	$\approx 227\,781\mu s$	$\approx 3\,851\mu s$
RyuJIT-x64	10^6	$\approx 47\,956\mu s$	$\approx 49\,670\mu s$	$\approx 3\,782\mu s$

Логичные вопросы:

- Почему PowerC работает так быстро?
- Почему PowerA работает очень быстро, но только на LegacyJIT-x86 при $N = 10^5$?

(Бенчмарки) Денормализованные числа

	N	PowerA	PowerB	PowerC
LegacyJIT-x86	10^5	$\approx 168\mu s$	$\approx 19\,264\mu s$	$\approx 370\mu s$
RyuJIT-x64	10^5	$\approx 3\,961\mu s$	$\approx 4\,142\mu s$	$\approx 370\mu s$
LegacyJIT-x86	10^6	$\approx 152\,770\mu s$	$\approx 227\,781\mu s$	$\approx 3\,851\mu s$
RyuJIT-x64	10^6	$\approx 47\,956\mu s$	$\approx 49\,670\mu s$	$\approx 3\,782\mu s$

Логичные вопросы:

- Почему PowerC работает так быстро?
- Почему PowerA работает очень быстро, но только на LegacyJIT-x86 при $N = 10^5$?

Рассмотрим два случая:

- LegacyJIT-x64, Mono-x64 (**SSE**); RyuJIT-x64 (**AVX**)
- LegacyJIT-x86 (**x87 FPU**)

(SSE/AVX) Денормализованные числа

i	PowerA	PowerC
0	1.0	1.0

3FF0000000000000000

3FF0000000000000000

(SSE/AVX) Денормализованные числа

i	PowerA	PowerC
0	1.0 3FF0000000000000	1.0 3FF0000000000000
1	0.96 3FEEB851EB851EB8	0.9600000000000008 3FEEB851EB851EB9

(SSE/AVX) Денормализованные числа

i	PowerA	PowerC
0	1.0 3FF0000000000000	1.0 3FF0000000000000
1	0.96 3FEEB851EB851EB8	0.9600000000000008 3FEEB851EB851EB9
885	2.0419318345555615E-16 3CAD6D6617566397	1.8041124150158794E-16 3CAA0000000000000
886	1.9602545611733389E-16 3CAC4010166769D8	1.6653345369377348E-16 3CA80000000000000
887	1.8818443787264053E-16 3CAB1EC7C3967A17	1.6653345369377348E-16 3CA80000000000000

(SSE/AVX) Денормализованные числа

i	PowerA	PowerC
0	1.0 3FF0000000000000	1.0 3FF0000000000000
1	0.96 3FEEB851EB851EB8	0.9600000000000008 3FEEB851EB851EB9
885	2.0419318345555615E-16 3CAD6D6617566397	1.8041124150158794E-16 3CAA0000000000000
886	1.9602545611733389E-16 3CAC4010166769D8	1.6653345369377348E-16 3CA80000000000000
887	1.8818443787264053E-16 3CAB1EC7C3967A17	1.6653345369377348E-16 3CA80000000000000
18171	6.42285339593621E-323 000000000000000D	1.6653345369377348E-16 3CA80000000000000
18172	5.92878775009496E-323 000000000000000C	1.6653345369377348E-16 3CA80000000000000
18173	5.92878775009496E-323 000000000000000C	1.6653345369377348E-16 3CA80000000000000

(x87 FPU) Денормализованные числа

```
; PowerA ( $10^5$  :  $\approx 167\mu s$        $10^6$  :  $\approx 151\,947\mu s$ )
fld    qword ptr ds:[14D2E28h] ; 0.96
fmulp st(1),st                ; Значение в регистре
; PowerB ( $10^5$  :  $\approx 19\,079\mu s$        $10^6$  :  $\approx 226\,219\mu s$ )
fld    qword ptr ds:[892E20h]   ; 0.96
fmul   qword ptr [ecx+4]
fstp   qword ptr [ecx+4]        ; Значение в памяти
```

(x87 FPU) Денормализованные числа

```
; PowerA ( $10^5 : \approx 167\mu s$             $10^6 : \approx 151\,947\mu s$ )
fld      qword ptr ds:[14D2E28h] ; 0.96
fmulp    st(1),st               ; Значение в регистре
; PowerB ( $10^5 : \approx 19\,079\mu s$         $10^6 : \approx 226\,219\mu s$ )
fld      qword ptr ds:[892E20h] ; 0.96
fmul    qword ptr [ecx+4]
fstp    qword ptr [ecx+4]        ; Значение в памяти
```

Intel® 64 and IA-32 Architectures Software Developer's Manual
§8.2 X87 FPU Data Types

With the exception of the 80-bit double extended-precision format, all of data types exist in memory only.

When they are loaded into x87 FPU data registers, they are converted into double extended-precision format and operated on in that format.

When a denormal number is used as a source operand, the x87 FPU automatically normalizes the number when it is converted to double extended-precision format.

Поговорим про суммирование

$$M = \sum_{n=0}^{\infty} 2^n \cdot \left[\frac{M_0}{2^n(n+1)} \right]$$



Сумма чисел

Задача. Дан `List<double>`. Необходимо реализовать функцию `Sum`, которая вернёт сумму элементов этого списка.

(Проблема) Сумма чисел

Имеется следующая проблема¹:

$$10^{16} + 1 = 10^{16}$$

¹Научное название: absorption

(Проблема) Сумма чисел

Имеется следующая проблема¹:

$$10^{16} + 1 = 10^{16}$$

Вспоминаем мат. часть:

	Digits	Lower	Upper
double	≈ 15.9	$2.3 \cdot 10^{-308}$	$1.7 \cdot 10^{+308}$

¹Научное название: absorption

(Алгоритм Кэхэна) Сумма чисел

```
double KahanSum(this List<double> list)
{
    double sum = 0, error = 0;
    foreach (var x in list)
    {
        var y = x - error;
        var newSum = sum + y;
        error = (newSum - sum) - y;
        sum = newSum;
    }
    return sum;
}
```

(Пример) Сумма чисел

```
var a = new List<double>()
{ 1016, 1, 1, ..., 1 };
    100 раз
WriteLine("{0:N}", a.Sum());
WriteLine("{0:N}", a.KahanSum());
```

// 10,000,000,000,000,000.00
// 10,000,000,000,000,100.00

(Пример) Сумма чисел

```
var a = new List<double>()
{ 1016, 1, 1, ..., 1 };
    100 раз
WriteLine("{0:N}", a.Sum());
WriteLine("{0:N}", a.KahanSum());

// 10,000,000,000,000,000.00
// 10,000,000,000,000,100.00
```

Проблема возникает в самых разных местах:

- Банковская сфера
- Страховая сфера
- Численные алгоритмы
- Ситуации с операциями над числами разных порядков

(А теперь про скорость) Сумма чисел

Задача. Дан `double[]` длины 4000. Необходимо реализовать функцию `Sum`, которая *быстро* вернёт сумму элементов массива.

(А теперь про скорость) Сумма чисел

Задача. Дан `double[]` длины 4000. Необходимо реализовать функцию `Sum`, которая *быстро* вернёт сумму элементов массива.

Решение А.

```
public double Bcl() => a.Sum();
```

(А теперь про скорость) Сумма чисел

Задача. Дан `double[]` длины 4000. Необходимо реализовать функцию `Sum`, которая *быстро* вернёт сумму элементов массива.

Решение А.

```
public double Bcl() => a.Sum();
```

Решение В.

```
public double Loop() {
    double sum = 0;
    for (int i = 0; i < a.Length; i++)
        sum += a[i];
    return sum;
}
```

(А теперь про скорость) Сумма чисел

Задача. Дан `double[]` длины 4000. Необходимо реализовать функцию `Sum`, которая *быстро* вернёт сумму элементов массива.

Решение А.

```
public double Bcl() => a.Sum();
```

Решение В.

```
public double Loop() {
    double sum = 0;
    for (int i = 0; i < a.Length; i++)
        sum += a[i];
    return sum;
}
```

Решение С.

```
public double Unroll() {
    double sum = 0;
    for (int i = 0; i < a.Length; i += 4)
        sum += a[i] + a[i + 1] + a[i + 2] + a[i + 3];
    return sum;
}
```

(Бенчмарки) Сумма чисел

```
const int N = 4000; double[] a; // random double[N] array
[Benchmark] double Bcl() => a.Sum();
[Benchmark] double Loop() {
    double sum = 0;
    for (int i = 0; i < a.Length; i++)
        sum += a[i];
    return sum;
}
[Benchmark] double Unroll() {
    double sum = 0;
    for (int i = 0; i < a.Length; i += 4)
        sum += a[i] + a[i + 1] + a[i + 2] + a[i + 3];
    return sum;
}
```

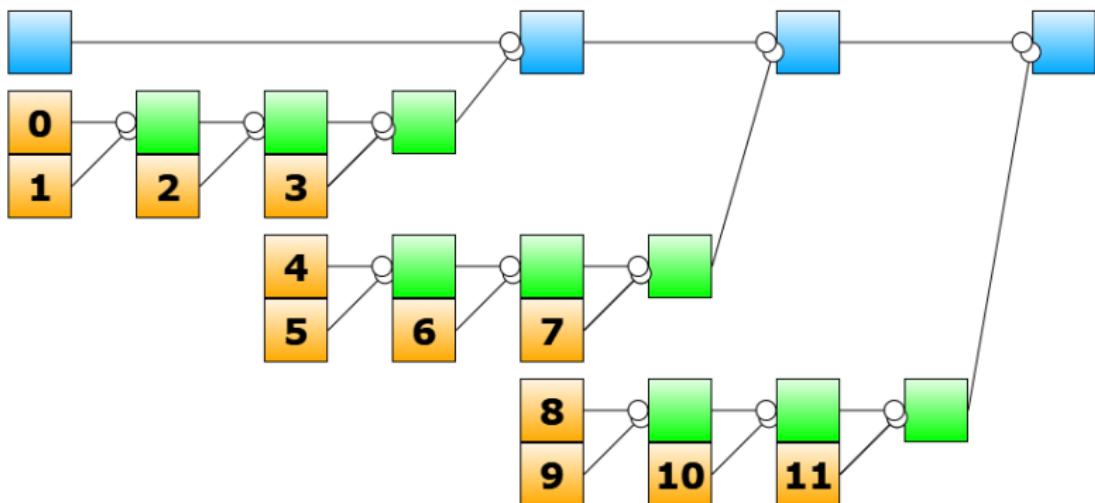
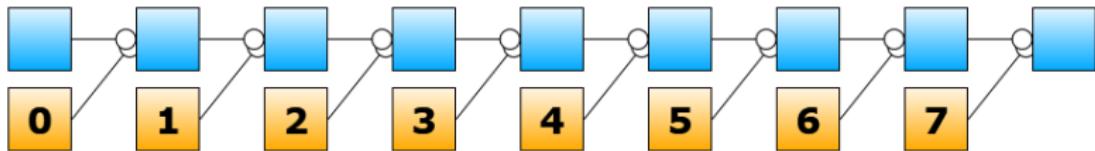
(Бенчмарки) Сумма чисел

```
const int N = 4000; double[] a; // random double[N] array
[Benchmark] double Bcl() => a.Sum();
[Benchmark] double Loop() {
    double sum = 0;
    for (int i = 0; i < a.Length; i++)
        sum += a[i];
    return sum;
}
[Benchmark] double Unroll() {
    double sum = 0;
    for (int i = 0; i < a.Length; i += 4)
        sum += a[i] + a[i + 1] + a[i + 2] + a[i + 3];
    return sum;
}
```

	Bcl	Loop	Unroll
LegacyJIT-x86	≈23.4μs	≈4.0μs	≈1.6μs
LegacyJIT-x64	≈28.3μs	≈4.0μs	≈1.8μs
RyuJIT-x64	≈31.1μs	≈4.2μs	≈2.4μs
Mono4.6.2-x86	≈72.7μs	≈14.3μs	≈2.2μs
Mono4.6.2-x64	≈54.7μs	≈12.3μs	≈3.6μs

(Примерная схема) Сумма чисел

Instruction Level Parallelism спешит на помощь!



Методическая литература

2³² + 1 = 641 6790412 $\text{sign}(x, y) \cdot ((x \& y) \cdot ((x \& y) \cdot \delta))$
2³² + 1 = 274777 678042110721
 $x - y \neq x + y + 1$
 $\lfloor a \rfloor + \lfloor b \rfloor \leq \lfloor a + b \rfloor \leq \lfloor a \rfloor + \lfloor b \rfloor + 1$ $\text{popcnt} = \sum_{i=0}^{n-1} \lfloor \frac{n}{2^i} \rfloor$
George Boole
(1815-1864) $\lfloor \sqrt{11111111} \rfloor = 1111$
 $(x \oplus 0) \oplus (x \cdot \delta) = x$ $\max(x, y, m) = ((x \oplus y) \& m) \oplus y$
 $d(n, d) = d(n-1, d-1), d \text{ even}$ $-x = x + 1$

Hacker's Delight

SECOND EDITION

Henry S. Warren Jr.

This standard specifies interchange and arithmetic formats and methods for binary and decimal floating-point arithmetic in computer programming environments. This standard specifies exception conditions and their default handling. An implementation of a floating-point system conforming to this standard may be realized entirely in software, entirely in hardware, or in any combination of software and hardware. For operations specified in the normative part of this standard, numerical results and exceptions are uniquely determined by the values of the input data, sequence of operations, and destination formats, all under user control.

Date of Publication : Aug. 29 2008
Status : Active
INSPEC Accession Number: 10369688
DOI: 10.1109/IEEESTD.2008.4610835

Download Chapters
Email
Print
Request Permissions
Report

f

Browse Standards - IEEE Std 754-2008

754-2008 - IEEE Standard for Floating-Point Arithmetic

Full Text
Sign-in or Purchase

Revision of ANSI/IEEE Std 754-1985

Status: Active - Approved

Abstract References Versions Definitions Cited By Keywords

This standard specifies interchange and arithmetic formats and methods for binary and decimal floating-point arithmetic in computer programming environments. This standard specifies exception conditions and their default handling. An implementation of a floating-point system conforming to this standard may be realized entirely in software, entirely in hardware, or in any combination of software and hardware. For operations specified in the normative part of this standard, numerical results and exceptions are uniquely determined by the values of the input data, sequence of operations, and destination formats, all under user control.

Date of Publication : Aug. 29 2008
Status : Active
INSPEC Accession Number: 10369688
DOI: 10.1109/IEEESTD.2008.4610835



Вопросы?

Андрей Акиньшин, JetBrains

<http://aakinshin.net>

<https://github.com/AndreyAkinshin>

https://twitter.com/andrey_akinshin

andrey.akinshin@gmail.com