

Министерство образования Республики Беларусь
Учреждение образования
БЕЛОРУССКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ
ИНФОРМАТИКИ И РАДИОЭЛЕКТРОНИКИ

Факультет компьютерных систем и сетей
Кафедра информатики
Дисциплина: Архитектура вычислительных систем

ОТЧЁТ
к практической работе
на тему

**АРИФМЕТИЧЕСКИЕ ОПЕРАЦИИ С ЧИСЛАМИ
С ПЛАВАЮЩЕЙ ТОЧКОЙ**

Выполнил: студент группы 153501
Тимофеев Кирилл Андреевич

Проверила:
Калиновская Анастасия Александровна

МИНСК 2023

СОДЕРЖАНИЕ

1. Краткие теоретические сведения
2. Задания и вывод программы
3. Программный код
4. Выводы

КРАТКИЕ ТЕОРЕТИЧЕСКИЕ СВЕДЕНИЯ

В формате с фиксированной точкой, в частности в дополнительном коде, можно представлять положительные и отрицательные числа в диапазоне, симметричном на числовой оси относительно точки 0. Расположив воображаемую Разделяющую точку в середине разрядной сетки, можно в этом формате представлять не только целые, но и смешанные числа, а также дроби.

Однако такой подход позволяет представить на ограниченной разрядной сетке множество вещественных чисел в довольно узком диапазоне. Нельзя представить очень большие числа или очень маленькие. При выполнении деления двух больших чисел, как правило, теряется дробная часть частного.

При работе в десятичной системе счисления ученые давно нашли выход из положения, применяя для представления числовых величин так называемую научную нотацию. Так, число 976 000000 000 000 можно представить в виде 9.76×10^{14} , а число 0,000000 000 000 0976 - в виде 9.76×10^{-14} . При этом, фактически, разделительная точка динамически сдвигается в удобное место, а для того чтобы "уследить" за ее положением в качестве второго множителя - характеристики, - используется степень числа 10 (основания характеристики). Это позволяет с помощью небольшого числа цифр (т.е. чисел с ограниченной разрядностью) с успехом представлять как очень большие, так и очень малые величины.

Этот же подход можно применить и в двоичной системе счисления. Число можно представить в виде

$$\pm S \times B^{\pm E}$$

Компоненты такого представления можно сохранить в двоичном слове, состоящем из трех полей:

- поле знака числа (плюс или минус)
- поле мантиссы S;
- поле порядка E;

Основание B подразумевается неявно и не сохраняется.

Принципы представления лучше пояснить на примерах. В крайнем левом бите слова хранится знак числа (0 – положительное, 1- отрицательное). В следующих 8 битах хранится значение порядка. Для представления порядка используется так называемый смещенный формат. Для получения действительного двоичного кода порядка из значения, сохраняемого в этом поле нужно вычесть фиксированное смещение. Как правило, смещение равно $(2^{k-1}-1)$, где k – разрядность поля порядка.. В данном случае $k = 8$, и в поле порядка можно представить коды в диапазоне от 0 до 255. Если принять значение смещения 127, то действительное значение порядка чисел, представленных в таком формате может находиться в интервале от -127 до +128. Стандарт IEEE формата с плавающей точкой

Для унификации формата представления чисел с плавающей точкой, что является необходимым условием переносимости программного обеспечения, Институтом инженеров по электротехнике и радиоэлектронике IEEE разработан стандарт 754 . В последнее десятилетие практически все процессоры и арифметические сопроцессоры проектируются с учетом требований этого стандарта.

Стандарт специфицирует два варианта формата: 32-битовый — обычной точности представления и 64-битовый — удвоенной точности представления. В первом формате поле порядка занимает 8 бит, а во втором -11 бит. Стандарт регламентирует использование числа 2 в качестве неявно заданного значения основания характеристики. Помимо основных, в стандарте предусмотрены два расширенных варианта форматов обычной и удвоенной точности, конкретная спецификация которых зависит от реализации вычислительной системы. Расширенные форматы позволяют включать дополнительные биты в поле порядка (расширение диапазона представления) в поле мантиссы (повышение точности представления). Расширенные форматы предназначены для промежуточных вычислений. За счёт повышения точности снижается вероятность появления ошибок округления, а при расширении диапазона снижается вероятность появления ошибки переполнения.

Таблица 3. Параметры форматов, регламентированные стандартом IEEE 754

Параметр	Формат			
	Обычная точность	Расширенный обычной точности	Удвоенная точность	Расширенный удвоенной точности
Размер слова (бит)	32	≥ 43	64	≥ 79
Поле порядка (бит)	8	≥ 11	11	≥ 15
Смещение порядка	127	Не регламентируется	1023	

Максимальное значение порядка	127	≥ 1023	1023	≥ 16383
Минимальное значение порядка	-126	≤ -1022	-1022	≤ -16382
Диапазон представления (по основанию 10)	$10^{-38}, 10^{+38}$	Не регламентируется	$10^{-308}, 10^{+308}$	Не регламентируется
Поле мантиссы (бит)	23	≥ 31	52	≥ 63
Количество значений порядка	254	Не регламентируется	2046	Не регламентируется
Количество значимый мантиссы	2^{23}	Не регламентируется	2^{52}	Не регламентируется
Количество отличающихся представимых величин	1.98×2^{31}	Не регламентируется	1.99×2^{63}	Не регламентируется

Сложение и вычитание

Алгоритмы выполнения операций сложения и вычитания в формате с плавающей точкой сложнее, чем аналогичные алгоритмы для чисел в формате с фиксированной точкой. Связано это, в первую очередь, с необходимостью выравнивания порядков операндов. Алгоритм включает четыре основных этапа.

1. Проверка на нуль.
2. Сдвиг мантисс для выравнивания порядков.
3. Суммирование или вычитание мантисс.
4. Нормализация результата.

Блок-схема типового алгоритма представлена на рис. 9. Детальный пошаговый анализ этого алгоритма покажет, какие функции используются при выполнении операций сложения и вычитания чисел в формате с плавающей точкой. В дальнейшем для определенности будем считать, что используется формат, регламентированный стандартом IEEE 754. Перед началом выполнения операций операнды должны быть помещены в регистры АЛУ. Если в используемом формате с плавающей точкой предполагается неявный старший разряд мантиссы, этот разряд должен быть в явном виде включен в регистры операндов, и все операции с ним в дальнейшем будут проводиться точно так же, как и с остальными разрядами мантиссы.

Поскольку операции сложения и вычитания отличаются только тем, что при вычитании предварительно изменяется знак второго операнда (вычитаемого), эта операция включена в ветвь ВЫЧИТАНИЕ на схеме алгоритма, после чего обе ветви сливаются. Далее анализируется, не равен ли один из операндов нулю. Если это так, то результат — значение второго операнда.

Следующий этап — изменение кодов операндов таким образом, чтобы значения их порядков стали равны. Чтобы понять, зачем это нужно, рас-

смотрим следующий пример в десятичной системе счисления:

$$(123 \times 10^0) + (456 \times 10^{-2}).$$

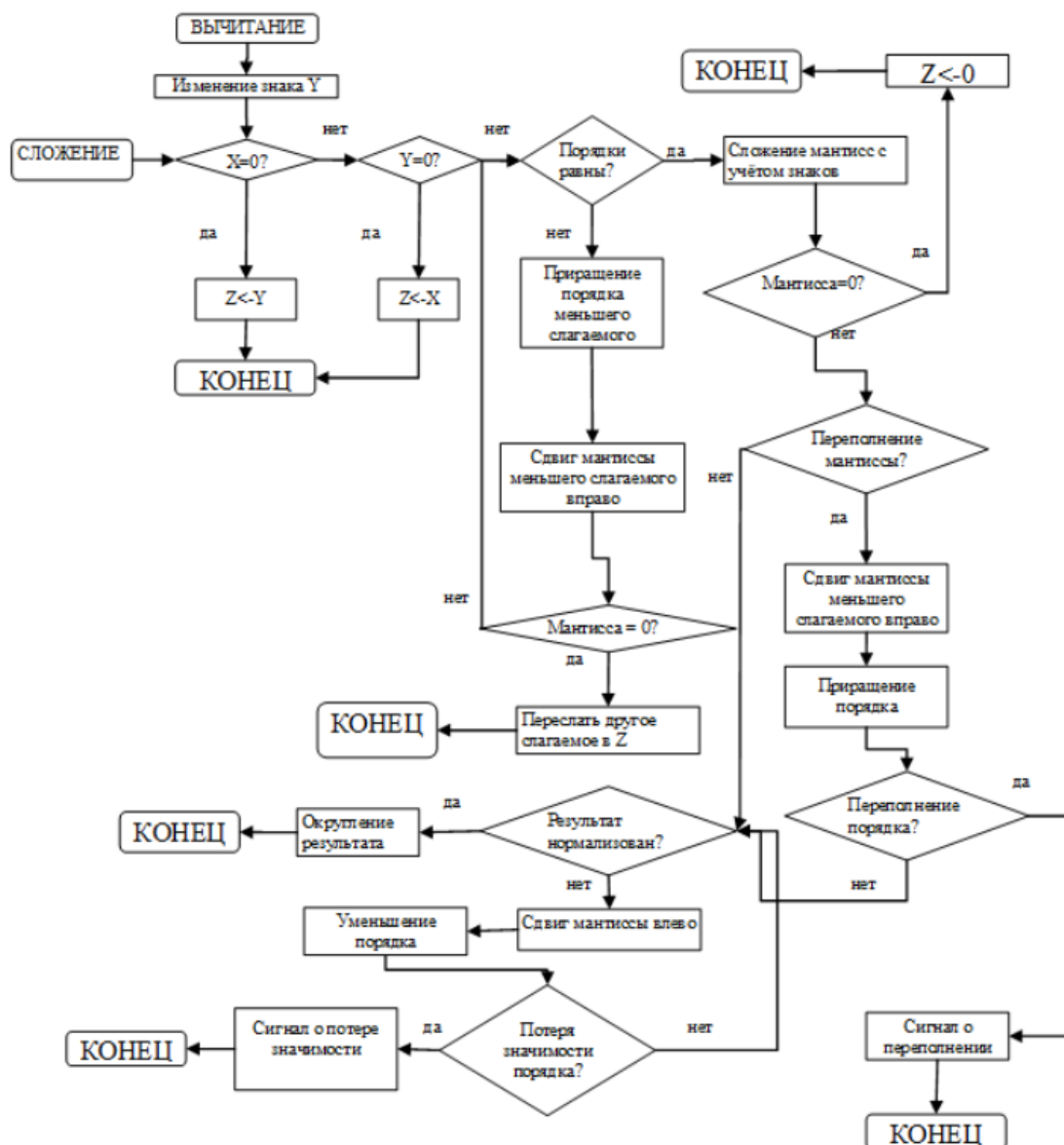
Очевидно, что нельзя просто сложить мантиссы этих двух чисел. Сначала нужно выровнять разрядные сетки обеих мантисс так, чтобы соответственные разряды (разряды с равным весом) занимали одинаковые позиции, т.е. цифра 4 второго числа находилась в той же позиции, что и цифра 3 первого числа. При этом порядки обоих чисел будут равны. Равенство порядков и есть, с точки зрения математики, условие, позволяющее складывать мантиссы обоих чисел в такой форме представления. Следовательно,

$$(123 \times 10^0) + (456 \times 10^{-2}) = (123 \times 10^0) + (4.56 \times 10^0) = 127.56 \times 10^0.$$

Выравнивание выполняется за счет сдвига мантиссы меньшего числа вправо или мантиссы большего числа — влево. Поскольку в любом варианте теряются цифры операнда, выполняется сдвиг мантиссы меньшего числа вправо, что приводит к утере ее младших разрядов. Одновременно со сдвигом мантиссы вправо порядок меньшего числа увеличивается. Сдвиги выполняются до тех пор пока значения порядков обоих чисел не станут равны.

После того, как порядки будут выровнены, наступает этап сложения мантисс с учетом их знаков. Поскольку слагаемые могут иметь разные знаки, их алгебраическая сумма может оказаться равной нулю. Не исключено и появление переполнения — переноса из старшего разряда суммы. В этом случае необходимо выполнить дополнительный сдвиг результата вправо и одновременно увеличить на 1 значение порядка. Но тогда возможно появление переполнения порядка, которое расценивается как аварийная ситуация и влечет за собой прекращение операции.

После сложения мантисс наступает этап нормализации результата. Нормализация представляет собой серию сдвигов кода мантиссы влево (в сторону старших разрядов) с одновременным уменьшением значения порядка до тех пор, пока значение старшей цифры мантиссы не станет отличным от нуля. Я специально оговариваю цифры, поскольку в случае, когда основанием характеристики является число 16, шестнадцатеричной цифрой мантиссы будет служить 4-разрядный двоичный код. Последняя операция — округление результата.



Умножение и деление

При работе с числами в формате с плавающей точкой алгоритмы умножения и деления оказываются проще алгоритмов сложения и вычитания.

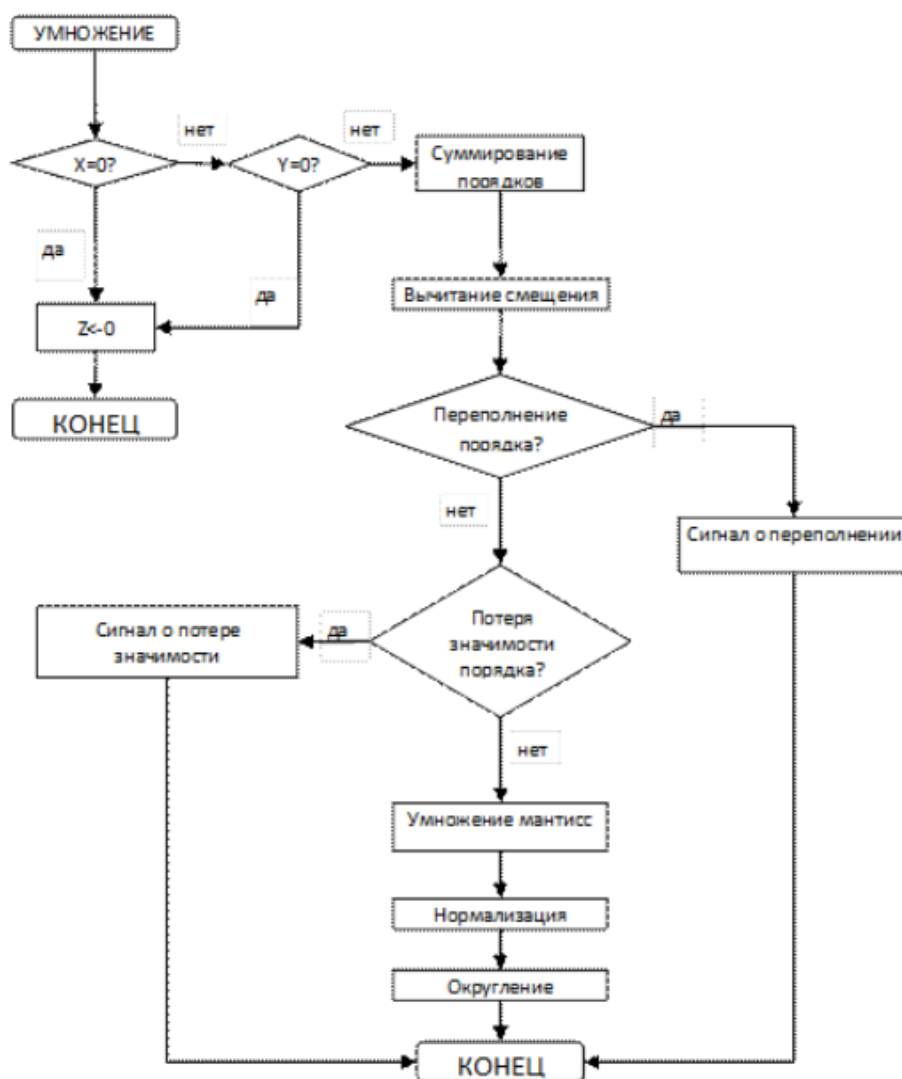


Рис. 10. Схема алгоритма умножения чисел в формате с плавающей точкой
($Z \leftarrow X \times Y$)

Сначала рассмотрим алгоритм умножения (рис. К.10). Сразу же после начала операции проверяется, не равен ли нулю один из сомножителей. Если это так, то произведение также будет равно нулю. Следующий шаг — суммирование порядков. Поскольку, как правило, для хранения порядков используется смещенное представление, при суммировании двух смещенных представлений результат будет смещен дважды. Поэтому после суммирования кодов порядков из суммы вычитается значение смещения. При суммировании может возникнуть как переполнение порядка, так и потеря значимости. В обоих случаях формируется соответствующий сигнал. Если порядок произведения не выходит из диапазона, определенного форматом, далее перемножаются мантиссы сомножителей с учетом их знаков. Умножение мантисс выполняется по тому же алгоритму, что и умножение целых чисел в прямом коде, т.е. фактически перемножаются числа без знака, а затем произведению приписывается знак "плюс" или "минус" в зависимости от сочетания знаков сомножителей. Произведение мантисс имеет разрядность, вдвое большую,

чем каждый из сомножителей. Лишние младшие разряды отбрасываются при округлении.

После того как будет получено произведение мантисс, результат нормализуется и округляется. Эти операции выполняются так же, как и при сложении или вычитании. Необходимо учесть, что при нормализации может возникнуть переполнение или потеря значимости порядка.

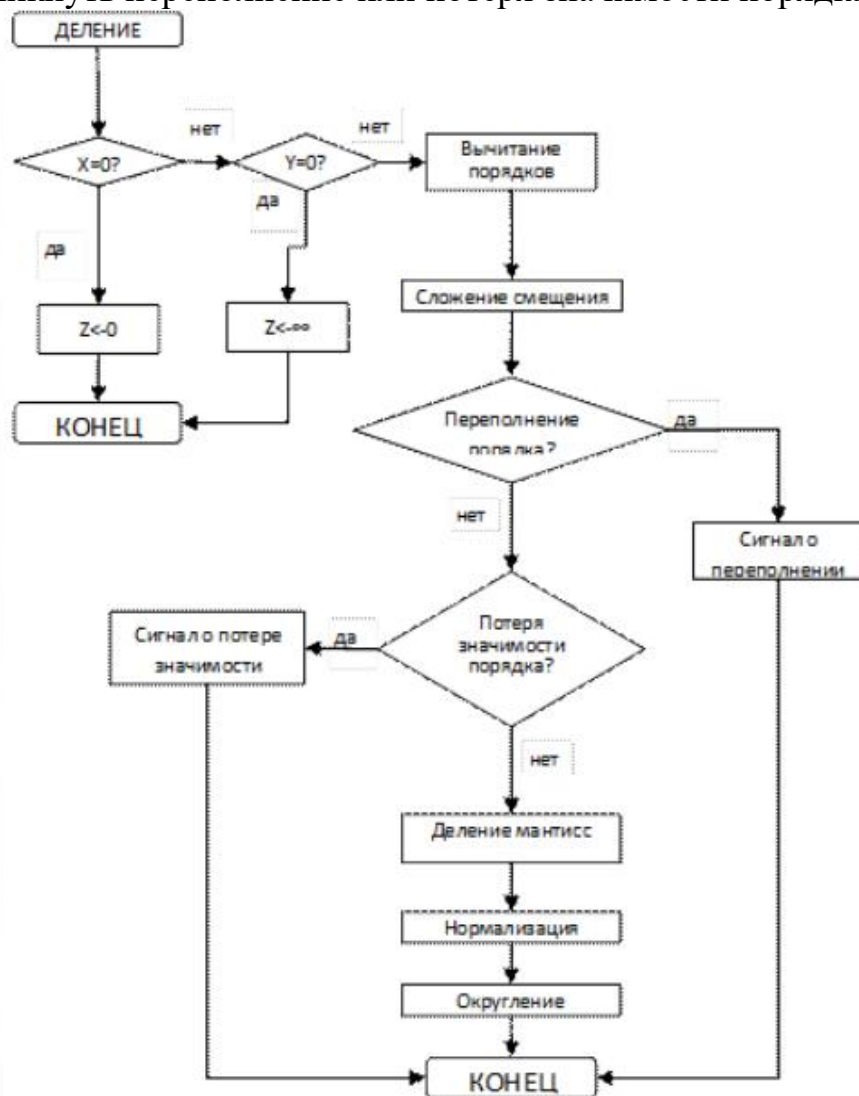


Рис. 11. Схема алгоритма деления чисел в формате с плавающей точкой ($Z \leftarrow X / Y$)

Теперь рассмотрим алгоритм деления (рис. К.11). Как и ранее, первый этап — анализ операндов на равенство нулю. Если нулю равно делимое, то результату сразу присваивается значение 0. Если же нулю равен делитель, то в зависимости от конкретной реализации АЛУ результату может быть присвоено значение "бесконечность" с соответствующим знаком или сформирован сигнал арифметической ошибки.

Следующий этап — вычитание кода порядка делителя из кода порядка делимого. При этом получится несмещенный код разности, который нужно скорректировать — сложить с кодом смещения. После завершения операций с порядком результата проверяется, не возникло ли переполнение порядка

или потеря значимости.

Следующий этап — деление мантисс. За ним следуют обычные операции нормализации и округления.

ЗАДАНИЯ И ВЫВОД ПРОГРАММЫ

Задание к лабораторной работе 4

Написать программу эмулятора АЛУ, реализующего *операции сложения и вычитания с плавающей точкой* над двумя введенными числами, с возможностью пошагового выполнения алгоритмов.

Задание к лабораторной работе 5

Написать программу эмулятора АЛУ, реализующего *операцию умножения с плавающей точкой* над двумя введенными числами, с возможностью пошагового выполнения алгоритмов.

Задание к лабораторной работе 6

Написать программу эмулятора АЛУ, реализующего *операцию деления с плавающей точкой* над двумя введенными числами, с возможностью пошагового выполнения алгоритмов.

Вывод программы

```

Введите первое число
44
Введите второе число
12
a = 0 10000010 100000000000000000000000
b = 0 10000010 100000000000000000000000
3145728
4194304

a+b
0 10000100 110000000000000000000000
1.7500000e5 = 56
3145728
4194304

a-b
0 00000000 000000000000000000000000
1.0000000e-127 = 0

a*b
0 10001000 000010000000000000000000
1.0312500e9 = 528

a/b
0 10000000 11010101010101010101010
1.8281250e1 = 3.66667

```

ПРОГРАММНЫЙ КОД

```

4 #include <iostream>
5 #include <bitset>
6 #include <string>
7
8 using namespace std;
9
10 const int SIZE = 32;
11 const int EN = 8;
12 const int MN = 23;
13 const int SHIFT = (1 << (EN - 1)) - 1;
14
15 typedef bitset<32> myfloat;
16 typedef bitset<EN> pow2;
17 typedef bitset<MN> mantis;
18
19 bitset<23> operator+(bitset<23> m1, bitset<23> m2)
20 {
21     return bitset<23>(m1.to_ulong() + m2.to_ulong());
22 }
23
24 bool isZero(myfloat a)
25 {
26     a[31] = 0;
27     return a.none();
28 }
29
30 myfloat zero(bool sign = 0)
31 {
32     myfloat rv = myfloat();
33     rv.reset();
34     rv[0] = sign;
35     return rv;
36 }
37
38 pow2 getPow(myfloat a)
39 {
40     pow2 rv = pow2();
41     for (int i = 0; i < EN; ++i)
42     {
43         rv[i] = a[i + MN];
44     }
45     return rv;
46 }
47
48 mantis getMantis(myfloat a)
49 {
50     mantis rv = mantis();
51     for (int i = 0; i < MN; ++i)
52     {
53         rv[i] = a[i];
54     }
55     return rv;
56 }
57
58 if (res >= 1 << 24)
59 {
60     res >>= 1;
61     pow_a += 1;
62     res -= 1 << 23;
63 }
64 myfloat rv = myfloat();
65 if (pow_a > 254)
66 {
67     cout << "Сигнал о переполнении";
68     return myfloat();
69 }
70 if (res == 0) return zero(sign);
71 setMantis(rv, res);
72 setPow(rv, pow_a);
73 rv[31] = sign;
74 return rv;
75 }
76
77 myfloat operator-(myfloat a, myfloat b)
78 {
79     b[31].flip();
80     return a + b;
81 }
82
83 myfloat operator*(myfloat a, myfloat b)
84 {
85     unsigned long pow_a = getPow(a).to_ulong();
86     unsigned long pow_b = getPow(b).to_ulong();
87     unsigned long pow_r = pow_a + pow_b - SHIFT;
88     unsigned long long mant_a = getMantis(a).to_ulong() + (1 << 23);
89     unsigned long long mant_b = getMantis(b).to_ulong() + (1 << 23);
90     unsigned long long mant_r = mant_a * mant_b;
91     mant_r >>= 23;
92     if (mant_r >= (1 << 24))
93     {
94         pow_r += 1;
95         mant_r >>= 1;
96     }
97     mant_r -= (1 << 23);
98     myfloat res = myfloat();
99     setMantis(res, mant_r);
100     setPow(res, pow_r);
101     res[31] = a[31] != b[31];
102     return res;
103 }
104
105 string toString(myfloat a)
106 {
107     long p = getPow(a).to_ulong() - 127;
108     mantis m = getMantis(a);
109     string str = (a[31] == 1) ? "-" : "";
110 }

```

```

60 void setMantis(myfloat& a, long m)
61 {
62     for (int i = 0; i < MN; ++i)
63     {
64         a[i] = m % 2;
65         m >>= 1;
66     }
67 }
68
69 void setPow(myfloat& a, unsigned long pow)
70 {
71     for (int i = 31 - EN; i < 31; ++i)
72     {
73         a[i] = pow % 2;
74         pow >>= 1;
75     }
76 }
77
78 void SetZero(myfloat& a)
79 {
80     a.reset();
81 }
82
83 myfloat operator+(myfloat a, myfloat b)
84 {
85     if (isZero(a)) return b;
86     if (isZero(b)) return a;
87     if (getPow(a).to_ulong() < getPow(b).to_ulong()) swap(a, b);
88     else if (getPow(a).to_ulong() == getPow(b).to_ulong())
89     {
90         if (getMantis(a).to_ulong() < getMantis(b).to_ulong()) swap(a, b);
91     }
92     unsigned long pow_a = getPow(a).to_ulong();
93     unsigned long pow_b = getPow(b).to_ulong();
94     cout << getMantis(a).to_ulong() << endl;
95     cout << getMantis(b).to_ulong() << endl;
96     long m_a = getMantis(a).to_ulong() + (1 << 23);
97     long m_b = getMantis(b).to_ulong() + (1 << 23);
98     while (pow_a != pow_b)
99     {
100         pow_b += 1;
101         m_b >>= 1;
102         if (m_b == 0) return a;
103     }
104     if (a[31] == 1) m_a ^= -1;
105     if (b[31] == 1) m_b ^= -1;
106     long res = m_a + m_b;
107     bool sign = res < 0;
108     if (sign) res = -res;
109     if (res >= 1 << 24)

```

```

170 string str = (a[31] == 1) ? "-" : "";
171 long long j = 5;
172 unsigned long long f = 1;
173 long long count = 1;
174
175 for (int i = m.size() - 1; i >= 16; --i)
176 {
177     f *= 10;
178     count *= 10;
179     if (m[i] == 1)
180     {
181         f += j;
182     }
183     j *= 5;
184 }
185
186 string mstr = to_string(f);
187 mstr.erase(mstr.begin());
188 return str + "1." + mstr + "e" + to_string(p);
189 }
190
191 string format(string bs)
192 {
193     bs.insert(1, 1, ' ');
194     bs.insert(2 + EN, 1, ' ');
195     return bs;
196 }
197
198 union fuck
199 {
200     float f;
201     unsigned long l;
202 };
203
204 myfloat operator/(myfloat a, myfloat b)
205 {
206     pow2 p_a = getPow(a);
207     pow2 p_b = getPow(b);
208     mantis m_a = getMantis(a);
209     mantis m_b = getMantis(b);
210     long res_pow = p_a.to_ulong() - p_b.to_ulong() + SHIFT;
211     unsigned long long ch = m_a.to_ulong() + (1 << 23);
212     unsigned long long zn = m_b.to_ulong() + (1 << 23);
213     ch <<= 23;
214     unsigned long long res_m = ch / zn;
215     if (!(res_m & (1 << 23)))
216     {
217         --res_pow;
218         res_m <<= 1;
219     }
220 }

```

```

224         res_m <<= 1;
225     }
226
227     res_m -= (1 << 23);
228     myfloat res = myfloat();
229     setMantis(res, res_m);
230     setPow(res, res_pow);
231
232     res[SIZE - 1] = a[31] != b[31];
233
234     return res;
235 }
236
237 int main()
238 {
239     setlocale(0, "");
240     fuck f;
241     cout << "Введите первое число\n";
242     cin >> f.f;
243     myfloat a = myfloat(f.f);
244     cout << "Введите второе число\n";
245     cin >> f.f;
246     myfloat b = myfloat(f.f);
247
248     cout << "a = " << format(bitset<32>(f.f).to_string()) << endl;
249     cout << "b = " << format(bitset<32>(f.f).to_string()) << endl;
250
251     auto r = a + b;
252     f.f = r.to_ulong();
253     cout << endl << endl << "a+b" << endl;
254     cout << format(r.to_string()) << endl;
255     cout << toString(r) << " " << f.f << endl;
256
257     r = a - b;
258     f.f = r.to_ulong();
259     cout << endl << endl << "a-b" << endl;
260     cout << format(r.to_string()) << endl;
261     cout << toString(r) << " " << f.f << endl;
262
263     r = a * b;
264     f.f = r.to_ulong();
265     cout << endl << endl << "a*b" << endl;
266     cout << format(r.to_string()) << endl;
267     cout << toString(r) << " " << f.f << endl;
268
269     r = a / b;
270     f.f = r.to_ulong();
271     cout << endl << endl << "a/b" << endl;
272     cout << format(r.to_string()) << endl;
273     cout << toString(r) << " " << f.f << endl;
274
275 }

```

ВЫВОДЫ

В результате данной работы были изучены методы представления чисел с плавающей точкой и формат IEEE-754.