

Возьмем несколько чисел и представим в формате с плавающей точкой. Несмотря, что для хранения мантииссы данного формата в стандарте IEEE754 используется 3 байта, ограничимся одним

A = -13,13; B = 46,46; C = -6,868; D = -0,4444; E = 0,07777 F = 0,1999.

Отдельно совершим переход целой части и после запятой.

1.

$$13_{10} = 1101_2;$$

Так уже 4 бита есть, то для оставшейся части, с учетом окружения, нужно определить 5 бит.

$$0,13 \times 2 = 0,26; 0,26 \times 2 = 0,52; 0,52 \times 2 = 1,04; 0,04 \times 2 = 0,08; 0,08 \times 2 = 0,16.$$

Поскольку при пятом умножении получился 0, то выделенные цифры формируют результат.  $0,13_{10} = 0010_2$ ;

Для формирования результата надо перенести запятую на три бита вперед. Тогда степень двойки будет 3. К степени следует добавить 127. и представить в двоичном коде.

$$130_{10} = 10000010_2$$

Окончательно получим

$$1 \text{ } 10000010 \text{ } 101001000...$$

**Замечание.** При компьютерном преобразовании будут, на самом деле, получены все 23 бита мантииссы.

2.

$$\begin{array}{r} 46 \overline{) 8} \\ 40 \overline{) 5} \\ \hline 6 \end{array}$$

$$46_{10} = 56_8 = 101110_2.$$

Еще нужно получить две цифры

$$0,46 \times 2 = 0,92; 0,92 \times 2 = 1,84; 0,84 \times 2 = 1,68.$$

С учетом ограничения получим

$$0,46_{10} = 011_2 = 10_2.$$

Для формирования результата надо перенести запятую на пять бит вперед. Тогда степень двойки будет 5. К степени следует добавить 127. и представить в двоичном коде.

$$132_{10} = 10000100_2$$

Окончательно получим

$$0 \text{ } 10000100 \text{ } 011101000...$$

3.

$$6_{10} = 110_2;$$

Так уже 3 бита есть, то для оставшейся части, с учетом окружения, нужно определить 6 бит.

$$0,868 \times 2 = 1,736; 0,736 \times 2 = 1,472; 0,472 \times 2 = 0,944; 0,944 \times 2 = 1,888; 0,888 \times 2 = 1,776; 0,776 \times 2 = 1,552.$$

Альтернативный вариант

$$6,944 \text{ } 7,552$$

Поскольку при пятом умножении получился 0, то выделенные цифры формируют результат.  $0,868_{10} = 110111_2 = 11100_2$ .

Альтернативный вариант

$$0,868 \times 8 = 6,944; 0,944 \times 8 = 7,552.$$

$$0,868_{10} = 67_8 = 110 \text{ } 111_2 = 11100_2$$

Для формирования результата надо перенести запятую на два бита вперед. Тогда степень двойки будет 2. К степени следует добавить 127. и представить в двоичном коде.

$$129_{10} = 1000001_2$$

Окончательно получим

$$1\ 1000001\ 101110000\dots$$

4.

У числа D целой части нет. Поэтому необходимо получить 8 бит из десятичной части числа

$$0,4444 \times 8 = 3,5552; 0,5552 \times 8 = 4,4416; 0,4416 \times 8 = 3,5328; 0,5328 \times 8 = 4,2624.$$

Тогда

$$0,4444_{10} = 3434_8 = 011\ 100\ 011\ 100_2.$$

Перенесем точку на один знак вправо. Степень получается -1. Добавим к степени 127.

$$126_{10} = 01111110_2.$$

Тогда окончательно мантисса после округления

$$11100011100_2 = 11100100_2.$$

Окончательно получим

$$1\ 01111110\ 110010000\dots$$

5.

У числа E целой части нет. Поэтому необходимо получить 8 бит из десятичной части числа

$$0,07777 \times 8 = 0,62216; 0,62216 \times 8 = 4,97728; 0,97728 \times 8 = 7,81824; 0,81824 \times 8 = 6,54592.$$

Тогда

$$0,07777_{10} = 0476_8 = 000\ 100\ 111\ 110_2.$$

Для нормализации перенесем точку на 4 бита вправо. Степень получается -4. Добавим к степени 127.

$$123_{10} = 01111011_2.$$

Тогда окончательно мантисса после округления

$$100111110_2 = 1010000_2.$$

Окончательно получим

$$0\ 01111011\ 101000000\dots$$

6.

У числа F целой части нет. Поэтому необходимо получить 8 бит из десятичной части числа

$$0,1999 \times 8 = 1,5992; 0,5992 \times 8 = 4,7936; 0,97728 \times 8 = 6,3488; 0,3488 \times 8 = 2,7904.$$

Тогда

$$0,1999_{10} = 1462_8 = 001\ 100\ 110\ 010_2.$$

Для нормализации перенесем точку на 2 бита вправо. Степень получается -2. Добавим к степени 127.

$$125_{10} = 01111101_2.$$

Тогда окончательно мантисса после округления

$$1100110010_2 = 11001101_2.$$

Окончательно получим

$$0\ 01111011\ 100110100\dots$$

### Сложение в формате с плавающей точкой

В отличие от сложения чисел с фиксированной точкой, когда может происходить только преобразование числа из отрицательного в положительное или наоборот для

выполнения нужной математической последовательности, то для чисел с плавающей точкой необходимо еще предварительно согласовать порядки, а после сложения проводить нормализацию результата.

Выполнил операцию сложения для чисел А, В и С в разной последовательности.

Представим числа в таблице, как это будет в задании. Значение мантиисы сокращены до шести бит с учетом округления.

	Знак	Порядок	Знак	Мантисса
А	0	0011	1	110101
В	0	0101	0	101111
С	0	0010	1	110111

Для сложения числа В и С следует уравнивать порядки путем сдвига мантиисы на разницу порядков. Сдвигать можно только мантиссу с меньшим порядком. Разрядная сетка должна сохраняться. После сдвига следует правильно округлить результат.

Поэтому С в модифицированном коде

$C = 11.110111_{\text{пк}} (0010 \text{ или } 2^2) = 11.00011011_{\text{пк}} (0101 \text{ или } 2^5) = 11.000111_{\text{пк}} (0101 \text{ или } 2^5)$ . К последнему биту был прибавлена 1 как следствие отброшенной единицы, которая была отброшена последней.

Перейдем в дополнительный код для выполнения операции сложения

$C = 11.000111_{\text{пк}} (0101 \text{ или } 2^5) = 11.111001_{\text{дк}} (0101 \text{ или } 2^5)$ . Сейчас уже можно провести сложение

	.	.	.	.	.	.	.	порядок
С		1	1	.	1	1	0 0 1	00.0101
В	+	0	0	.	1	1	0 1 1 1	00.0101
С+В		1	1	.	1	1	0 0 0 0	00.0101

Поскольку после битов знака положительного числа наблюдается 1, то число нормализовано и можно переходить к следующей операции.

Число А имеет меньший на 2 порядок, поэтому следует произвести сдвиг его мантиисы на 2.

$A = 11.110101_{\text{пк}} (0011 \text{ или } 2^3) = 11.00110101_{\text{пк}} (0101 \text{ или } 2^5)$  последним отброшенным является 0, то прибавлять 1 к последнему биту не нужно.

Перейдем в дополнительный код для выполнения операции сложения

$A = 11.001101_{\text{пк}} (0101 \text{ или } 2^5) = 11.110011_{\text{дк}} (0101 \text{ или } 2^5)$ . Сейчас уже можно провести сложение

	.	.	.	.	.	.	.	порядок
С+В		0	0	.	1	1	0 0 0 0	00.0101
А	+	1	1	.	1	1	0 0 1 1	00.0101
С+В+А		1	1	.	1	1	0 0 1 1	00.0101

Поскольку после битов знака положительного числа наблюдается 1, то число нормализовано.

Для перевода в стандарт IEEE754 прибавим к порядку 127

	.	.	.	.	.	.	.	.
5		0	0	.	0	0	0 0 0 1 0 1	
127	+	0	0	.	0	1	1 1 1 1 1 1	
		0	0	.	1	0	0 0 0 1 0 0	

Число положительное

$S+B+A = 0 10000100 \mathbf{000110000}...$

Выполним сложение в другой последовательности. Сложим сначала А и С.

Степень числа А больше на один, чем С. Поэтому сдвинем мантиссу числа С на один.

$C = 11.110111_{\text{пк}} (0010 \text{ или } 2^2) = 11.0110114_{\text{пк}} (0011 \text{ или } 2^3) = 11.011100_{\text{пк}} (0011 \text{ или } 2^3)$  К последнему биту был прибавлена 1 как следствие отброшенной единицы, которая была отброшена последней.

$C = 11.011100_{\text{пк}} (0011 \text{ или } 2^3) = 11.100100_{\text{дк}} (0011 \text{ или } 2^3)$ .

Число отрицательное. Преобразуем его в дополнительный код, чтобы можно было выполнить сложение

$A = 11.110101_{\text{пк}} (0011 \text{ или } 2^3) = 11.001011_{\text{дк}} (0011 \text{ или } 2^3)$

	.	.	.		порядок
C		1	1	. 1 0 0 1 0 0	00.0101
A	+	1	1	. 0 0 1 0 1 1	00.0101
C+A		1	0	. 1 0 1 1 1 1	00.0101

Число получилось с переполнением, поскольку наблюдается 10 в модифицированном коде знака. Требуется нормализация. Для нормализации следует сделать сдвиг и, как следствие увеличить степень на 1

$C+A = 10.101111_{\text{дк}} (0011 \text{ или } 2^3) = 11.010111_{\text{дк}} (0100 \text{ или } 2^4)$ .

Для сложения с числом В необходимо уравнивать порядки, и, в частности, порядок промежуточной суммы сдвинуть до степени 5, которая есть у В.

$C+A = 11.010111_{\text{дк}} (0100 \text{ или } 2^4) = 11.101011_{\text{дк}} (0101 \text{ или } 2^5)$ .

	.	.	.	.	.	.	.	.	порядок
C+A		1	1	. 1 0 1 0 1 1	00.0101				
B	+	0	0	. 1 1 0 1 1 1	00.0101				
C+B+A		1	0	. 1 0 0 0 1 0	00.0101				

Результат имеет отличие, что присуще операции сложения в формате с плавающей точкой

Проведем сложение для другой тройки чисел

	Знак	Порядок	Знак	Мантисса
D	1	0001	1	<b>111001</b>
E	1	0100	0	<b>101000</b>
F	1	0010	0	<b>110011</b>

Для сложения D и E надо сравнить порядки и сдвинуть мантиссу числа с меньшим порядком, т.е. E.

Число E имеет меньший на 3 порядок, поэтому следует произвести сдвиг его мантиссы на 3.

$E = 00.101000 (0101 \text{ или } 2^{-4}) = 00.000101000 (0001 \text{ или } 2^{-1})$  последним отброшенным является 0, то прибавлять 1 к последнему биту не нужно.

Для проведения сложения, переведем число D дополнительный код

$D = 11.111001_{\text{пк}} (0001 \text{ или } 2^{-1}) = 11.000111_{\text{пк}} (0001 \text{ или } 2^{-1})$

	.	.		порядок
E		0	0 . 0 0 0 1 0 1	11.0001 <sub>пк</sub>
D	+	1	1 . 0 0 0 1 1 1	11.0001 <sub>пк</sub>
D+E		1	1 . 0 0 1 1 0 0	11.0001 <sub>пк</sub>

Число нормализовано, так у результата отрицательного числа в дополнительном коде стоит 0 после модифицированного кода знака, поэтому нормализация не требуется.

Число F имеет меньший порядок чем итоговая сумма на один, поэтому сдвинем его мантиссу на 1.

$F = 00.110011 (0010 \text{ или } 2^{-2}) = 00.00110014 (0001 \text{ или } 2^{-1}) = 00.011010 (0001 \text{ или } 2^{-1})$  последним отброшенным является 1, то нужно прибавить 1 к последнему биту.

		.	.	.					порядок	
D+E	1	1	.	0	0	1	1	0	0	11.0001 <sub>пк</sub>
F	0	0	.	0	1	1	0	1	0	11.0001 <sub>пк</sub>
D+E+F	1	1	.	1	0	0	1	1	0	11.0001 <sub>пк</sub>

Число не нормализовано, так у результата отрицательного числа в дополнительном коде стоит 1 после модифицированного кода знака, поэтому произведем сдвиг и скорректируем порядок.

$$F = 11.100110_{\text{дк}} (0001 \text{ или } 2^{-1}) = 11.001100_{\text{дк}} (0010 \text{ или } 2^{-2})$$

Перейдем в прямой код

$$F = 11.001100_{\text{дк}} (0010 \text{ или } 2^{-2}) = 11.110100_{\text{пк}} (0010 \text{ или } 2^{-2})$$

Прибавим к порядку 127. для этого перейдем в обратный код (можно и в дополнительном) значение порядка при восьмибитовом представлении значения).

$$f_m = 11.00000010_{\text{пк}} = 11.11111101_{\text{ок}}$$

		.	.	.	.	.	.	.	.	
-2		1	1	.	1	1	1	1	0	1
127	+	0	0	.	0	1	1	1	1	1
		1	0	.	0	1	1	1	1	0

Так как это обратный код, то необходимо прибавить единицу перед модифицированным кодом знака к младшему биту результата.

-2		0	0	.	0	1	1	1	1	0	0
127	+	0	0	.	0	0	0	0	0	0	1
		0	0	.	0	1	1	1	1	0	1

$$D+E+F = 1\ 01111101\ \mathbf{1010000}\dots$$

	Знак	Порядок	Знак	Мантисса
A	0	0011	1	110101
B	0	0101	0	101111
D	1	0001	1	111001
E	1	0100	0	101000

Произведем умножение числа A на D, представленных в формате с плавающей точкой. Для процесса умножения существуют разные подходы. Используем алгоритм со сдвигом промежуточного результата. Знак результата можно определить сразу через операцию XOR знаков двух чисел. Поскольку оба числа отрицательные то произведение их будет положительное. По этой причине можно производить отдельное умножение мантисс без учета знака. Отдельно производится операция сложения порядков. Результат сложения может быть подкорректирован, если по результату перемножения мантисс потребуется нормализация.

.	00.110101	Модуль мантиссы числа A
	00.111001	Модуль мантиссы числа D
	00.000000	Заполнение регистра хранения промежуточного результата
+	00.110101	Результат умножения последнего бита числа D на число A
	00.110101	Промежуточная сумма
	00.011010	Результат сдвига влево на один
+	00.000000	Результат умножения второго бита с конца числа D на число A
	00.011010	Промежуточная сумма

+ 00.001101	0	Результат сдвига влево на один
00.000000		Результат умножения третьего бита с конца числа D на число A
00.001101		Промежуточная сумма
+ 00.000110	1	Результат сдвига влево на один
00.110101		Результат умножения четвертого бита с конца числа D на число A
00.111011		Промежуточная сумма
+ 00.011101	1	Результат сдвига влево на один
00.110101		Результат умножения пятого бита с конца числа D на число A
01.010010		Промежуточная сумма
00.101001	0	Результат сдвига влево на один
00.110101		Результат умножения последнего бита с конца числа D на число A
01.011110		Окончание процесса умножения (два бита знака указывают на переполнение)
00.101111	0	Результат нормализации (сдвига влево на один)

В буфере хранения последнего отброшенного бита 0, округление не изменит результат мантиссы. Однако степень следует увеличить на 1 из-за сдвига при нормализации. Произведем сложение порядков  $a_m$  и  $d_m$  чисел A и D для формирования результата. Порядок числа A три, а D минус два. Представим числа в обратном коде, при восьми битах значения числа

$$d_m = 11.00000010_{\text{пк}} = 11.11111101_{\text{ок}}$$

$$\begin{array}{r}
 \begin{array}{cccccccccccc}
 & \cdot & \cdot & & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 d_m & + & 1 & 1 & \cdot & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\
 a_m & & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\
 \hline
 & 1 & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0
 \end{array}
 \end{array}$$

По правилам обратного кода первую единицу, которая вышла за два бита знака следует прибавить к младшему биту

$$\begin{array}{r}
 \begin{array}{cccccccccccc}
 & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 + & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 \hline
 & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
 \end{array}
 \end{array}$$

К результату следует добавить еще единицу как результат нормализации мантиссы

$$\begin{array}{r}
 \begin{array}{cccccccccccc}
 & & & & & & & & & & \cdot & \\
 + & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 \hline
 & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0
 \end{array}
 \end{array}$$

И для представления в формате IEEE754 прибавим к последней сумме 127.

$$\begin{array}{r}
 \begin{array}{cccccccccccc}
 & & & & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\
 127 & + & 0 & 0 & \cdot & 0 & 1 & 1 & 1 & 1 & 1 & 1 \\
 2 & & 0 & 0 & \cdot & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\
 \hline
 & 0 & 0 & \cdot & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1
 \end{array}
 \end{array}$$

Тогда результат (первый бит указывает на знак, зеленные цифры это порядок + 127, первые 5 черных бит это мантисса произведения без первой единицы)

$$A \otimes D = 0 \text{ 10000001 01111000...}$$