



# ОПЕРАЦИИ НАД ЧИСЛАМИ С ПЛАВАЮЩЕЙ ТОЧКОЙ

Представление вещественных чисел

# Представление вещественных чисел

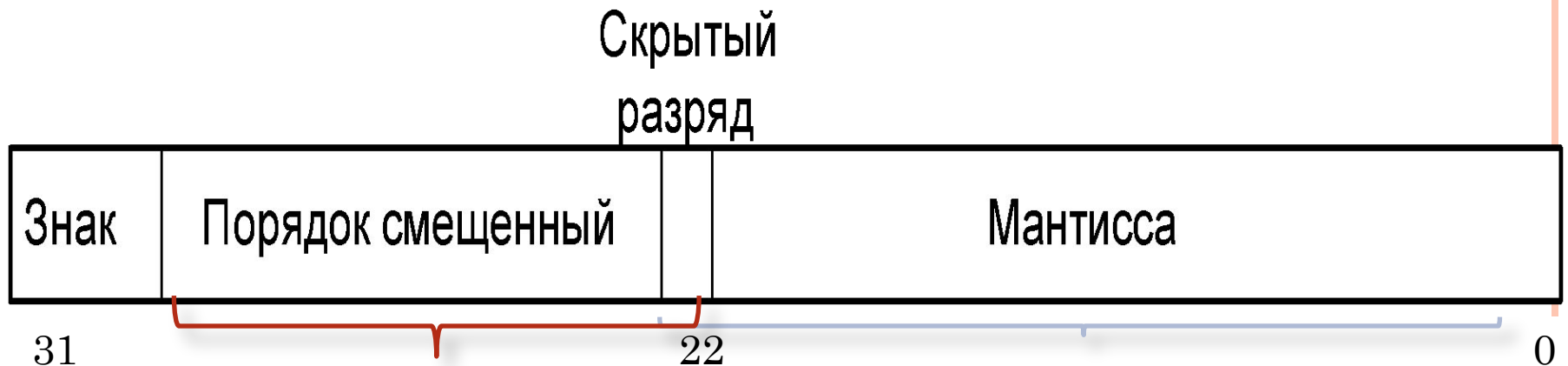
IEEE 754 (IEC 60559) — широко используемый стандарт IEEE, описывающий формат представления чисел с плавающей точкой. Используется в программных (компиляторы с разных языков программирования) и аппаратных (CPU и FPU) реализациях арифметических действий (математических операций).

$$X = \pm S^{P_x} q_x,$$

- *определяются требования к порядку- смещенный*  $P_{\text{смещ.}} = (2^{k-1}-1) + P_{\text{исх.}}$

*k*- кол-во разрядов выделенных под порядок

- *определяются требования к мантиссе-нормализованная*  $2 > |q_x| \geq 1$



# ПРЕДСТАВЛЕНИЕ ВЕЩЕСТВЕННЫХ ЧИСЕЛ

Базируется на экспоненциальной форме записи числа:

$$A = m * s^p$$

*m* - мантисса числа  
*s* - основание СС  
*p* - порядок числа

$$100_{10\text{с.с.}} = 0.1 * 10^3 = 10000 * 10^{-2}$$

1. определяются требования к порядку: может быть как + , так и -

2. определяются требования к мантиссе.

Для единообразия представления чисел используется **нормализованная форма**:

$1/s \leq |m| < 1$  (правильная дробь и после запятой цифра, отличная от нуля.)

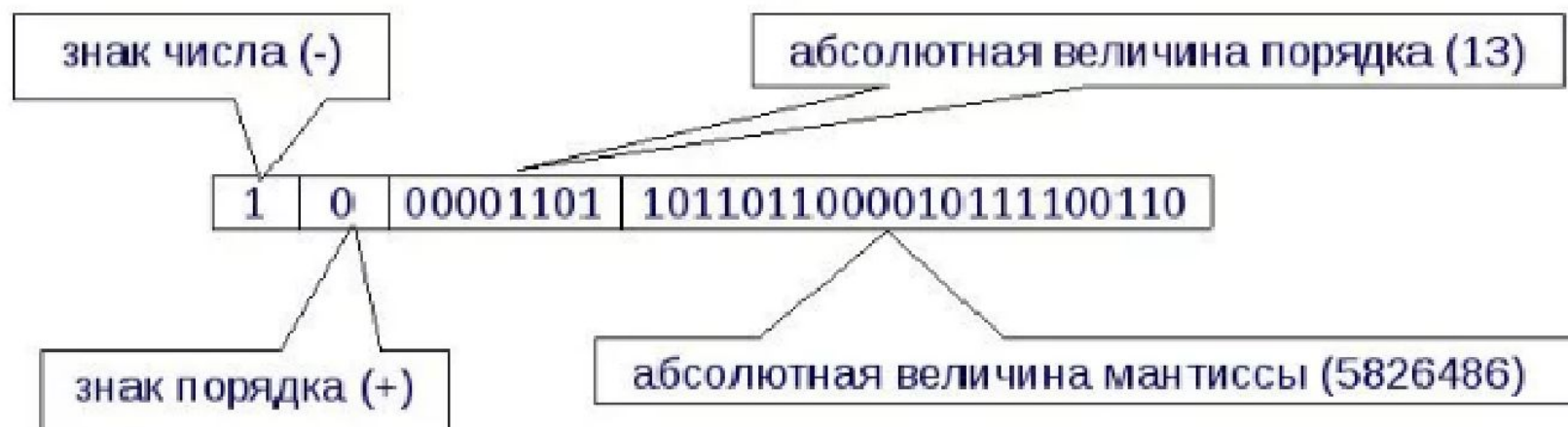
**Пример.** Преобразуйте число **555,55**, записанное в естественной форме, в экспоненциальную форму с нормализованной мантиссой: **555,55 = 0,55555 \* 10<sup>3</sup>**

Нормализованная мантисса: 0,55555

Порядок:  $p = 3$

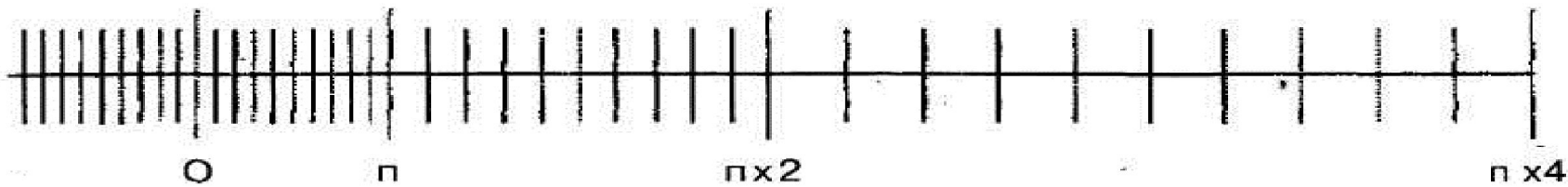


При представлении чисел с плавающей точкой в разрядах ячейки отводится место для знака числа, знака порядка, абсолютной величины порядка, абсолютной величины мантиссы.



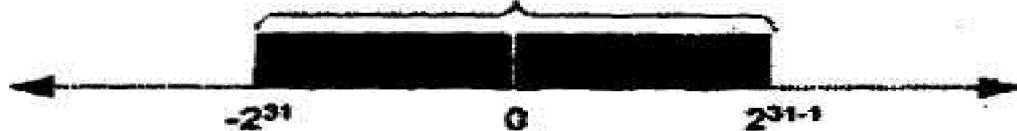
ячейке записано отрицательное двоичное число  $-1011011000010.11110011$   
десятичном представлении это будет число  $-5826.486$





Плотность чисел с плавающей запятой на числовой оси

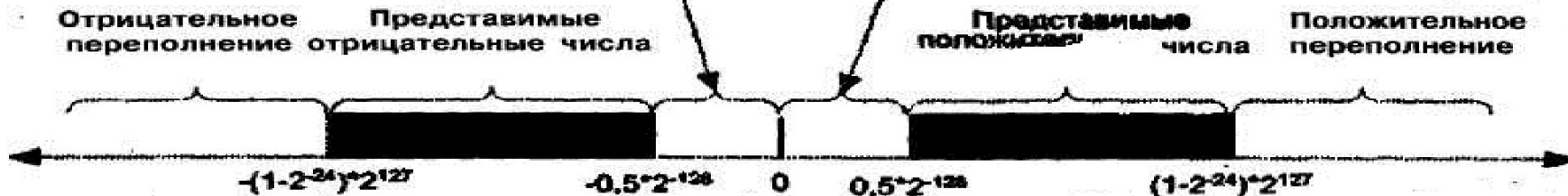
Представимые целые числа



С фиксированной точкой

Отрицательная потеря значимости

Положительная потеря значимости

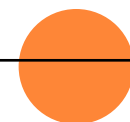


С плавающей точкой



В процессорах Intel (стандарт IEEE 754) применяется 3 (основных) формата с плавающей точкой: короткий, длинный и расширенный.

	Короткий	Длинный	Расширенный
Разрядность числа	32	64	80
Разрядность мантиссы	24	53	64
Диапазон значений	$10^{-38} \dots 10^{+38}$	$10^{-308} \dots 10^{+308}$	$10^{-4932} \dots 10^{+4932}$
Размерность порядка	8	11	15
Значение ( $2^{k-1}-1$ )	+127	+1023	+16383
Диапазон $p_{\text{смещ}}$	0...255	0...2047	0...32767
Диапазон порядков	-126...+127	-1022+1023	-16382+16383



## замечания

- *Рсм*, содержащий во всех разрядах 1 не используется, т.к. зарезервирован для указания на переполнение порядка или потерю значимости мантиссы.
- При этом + и - переполнение идентифицируются соответственно + и -мантиссами, содержащими в обоих случаях 1 во всех цифровых разрядах. Указанием на потерю значимости служит отрицательность порядка мантиссы с 0 во всех цифровых разрядах.
- При нулевых кодах порядка и мантиссы представляемое число полагается равным 0.
- Нулевому порядку в коротком формате соответствует значение фиксированного смещения равное 127, а в разрядной сетке запишется двоичное представление 01111111
- Отрицательному порядку -1, соответствует  $127-1=126 \rightarrow 01111110$  Положительному порядку +1, соответствует  $127+1=128 \rightarrow 10000000$ , т.е. все положительные порядки имеют старший бит порядка равный 1, а отрицательные-0.



Так как нормализованное число в старшем разряде всегда содержит 1, то при его представлении в памяти появляется возможность считать 1-й разряд вещественного числа единичным по умолчанию. И учитывать его наличие только на аппаратном уровне. Это дает возможность увеличить диапазон представимых чисел. Это утверждение справедливо только для короткого и длинного форматов.





$45.56_{10} = \mathbf{1011\ 01.10\ 0011\ 1101\ 0111\ 0000\ 1010\ 0011\ 1101\ 0111\ 0000\ 1010\ 1110\ 0001\ 0111\ 0001}$

Смещенный порядок  $127+5=132$  или  $10000100$

Запись числа в коротком формате со скрытым старшим разрядом мантииссы.

$0\ 10000100\ 01101100011110101110001$

Запись результата в 16-ричной с.с. **42363D71 H**



**ПРИМЕР** . ЗАПИСАТЬ ПРЕДСТАВЛЕНИЕ ЧИСЛА 0,089 В ФОРМЕ С ПЛАВАЮЩЕЙ ТОЧКОЙ.

### РЕШЕНИЕ

1. Переведем в двоичную СС:

$$0,089_{10} = 0,0001011000000000...0_2$$

2. Запишем в форме нормализованного двоичного числа:  $0,10110000000000000000 * 10_2^{100}$

3. Вычислим машинный порядок в двоичной СС:

$$P_{CM} = 01111111 - 100 = 01111011$$

3. Запишем число в коротком формате:

$$0 \ 01111011 \ 01100000000000000000$$

**Шестнадцатеричная форма 3DB00000**



# АЛГОРИТМ ЗАПИСИ ВНУТРЕННЕГО ПРЕДСТАВЛЕНИЯ ВЕЩЕСТВЕННОГО ЧИСЛА

1. Перевести модуль числа в двоичную СС с 24 значащими цифрами.
2. Нормализовать двоичное число.
3. Найти машинный порядок в двоичной СС.
4. Учитывая знак числа, записать его в 4-х байтовом машинном слове.



**ПРИМЕР 1** . ЗАПИСАТЬ ПРЕДСТАВЛЕНИЕ ЧИСЛА 122,1875 В ФОРМЕ С ПЛАВАЮЩЕЙ ТОЧКОЙ В КОРОТКОМ ФОРМАТЕ

1. Переведем в двоичную СС:

$$122,1875_{10} = 1111010,0011000000000000_2$$

2. Запишем в форме нормализованного двоичного числа:  $0,111101000110000000000000 * 10^6$

3. Вычислим  $R_{см}$  :

$$R_{см} = 127 + 6 \rightarrow 10000101_{2С.С}$$

3. Запишем число в 4-х байтовой ячейке:

0    10000101    111010    00110000    00000000

Шестнадцатеричная форма 42F46000





# ПРЕДСТАВЛЕНИЕ ЧИСЕЛ В ФОРМАТЕ С ПЛАВАЮЩЕЙ ТОЧКОЙ

Занимает в памяти ПК **4** (обычная точность) или **8 байтов** (*двойная точность*)

Выделяются разряды для хранения знака мантиссы, знака порядка, порядка и мантиссы.

Максимальное значение порядка числа:

$$1111111^2 = 127_{10}$$

Максимальное значение числа составляет:

$$2^{127} = 1,7014118346046923173168730371588 * 10^{38}$$

Максимальное значение положительной мантиссы равно:

$$2^{23} - 1 \sim 2^{23} = 2^{(10*2,3)} \sim 1000^{2,3} = 10^{(2,3*3)} \sim 10^7$$

Максимальное значение чисел обычной точности вычислений составляет  
 $1,701411 * 10^{38}$



## **Операции сложения над числами с плавающей точкой:**

$$Z=X+Y= S^{P_x} (q_x+q_y / S^{P_x-P_y}) = S^{P_z} q_z$$

### **Алгоритм сложения чисел с плавающей точкой:**

1. Производится выравнивание порядков чисел. Порядок меньшего (по модулю) числа принимается равным порядку большего, а мантисса меньшего сдвигается вправо на количество разрядов, равное разности порядков.
2. Производится сложение (вычитание) мантисс по правилам двоичной арифметики.
3. Нормализация результата.



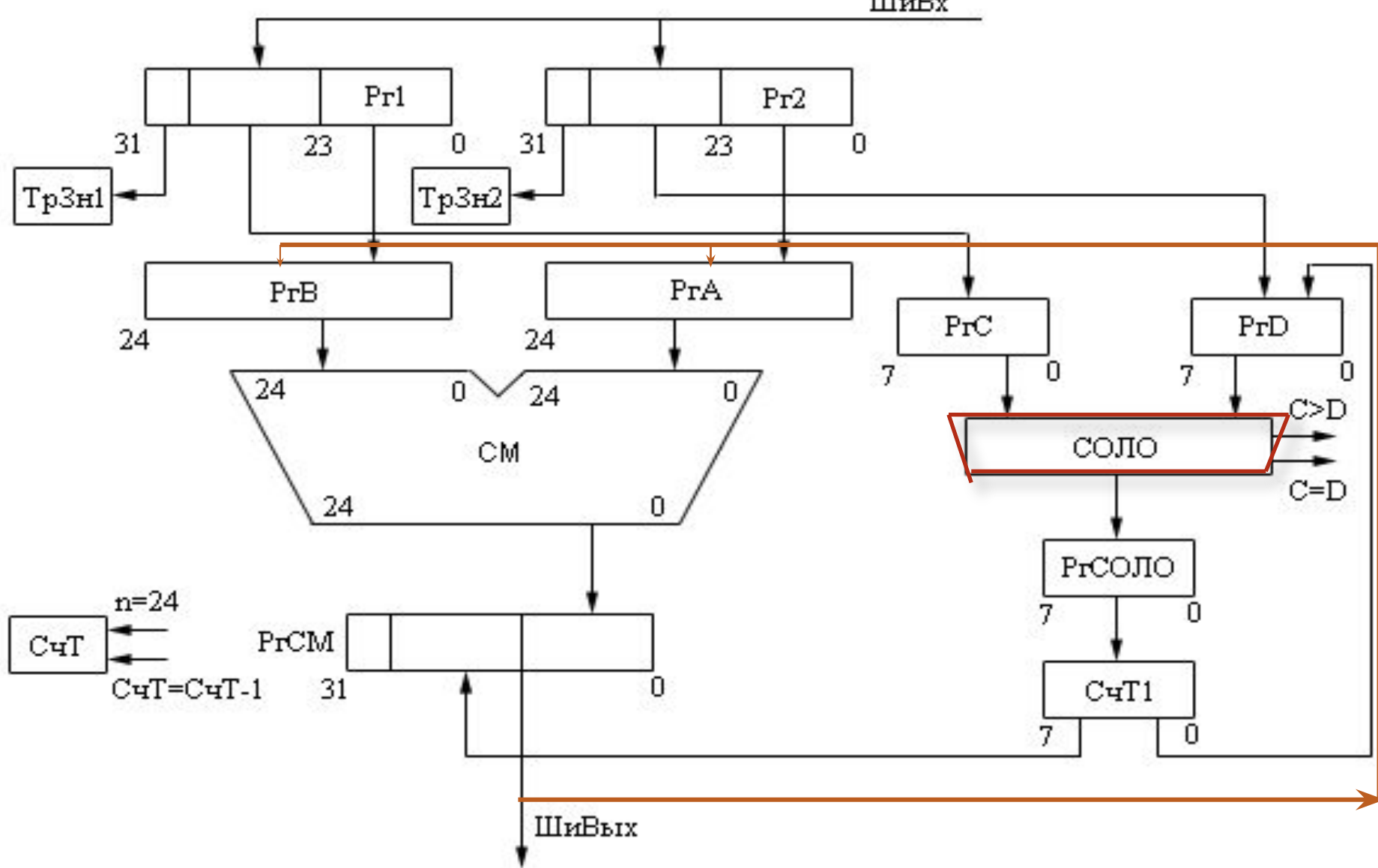
При сравнении порядков возможны пять случаев:

1.  $P_x - P_y > m$  ( $m$  — число разрядов мантиссы). В качестве результата суммирования сразу может быть взято первое слагаемое, так как при выравнивании порядков все разряды мантиссы второго слагаемого принимают нулевое значение;
2.  $P_y - P_x > m$  - В качестве результата суммирования может быть взято второе слагаемое;
3.  $P_x - P_y = 0$ . Можно приступить к суммированию мантисс;
4.  $P_x - P_y = k_1 (k_1 < m)$ . Мантисса второго слагаемого сдвигается на  $k_1$  разрядов вправо, затем производится суммирование мантисс;
5.  $P_y - P_x > = k_2 (k_2 < m)$ . Перед выполнением суммирования мантисс производится сдвиг на  $k_2$  разрядов вправо мантиссы первого слагаемого.

Сложение (вычитание) мантисс производится по правилам сложения (вычитания) чисел с фиксированной точкой.







Нормализация суммы (разности) производится в случае невыполнения условия  $1 > qz \geq 1/s$ , при этом,

-если  $qz \geq X$ ,  $Pz$  увеличивается на 1, а мантисса  $qz$  сдвигается на один  $S$ -ичный разряд вправо, что дает  $|qz| < 1$ .

-если  $|qz| < 1$ , то мантисса результата сдвигается на разряд влево при одновременном уменьшении порядка результата на 1. Эти операции производятся до тех пор, пока не станет выполняться условие  $qz \geq 1/s$ .

(При  $qz = 0$  нормализация не выполняется.)

При получении порядка  $+pz$ , переполняющего разрядную сетку, должен формироваться сигнал прерывания из-за переполнения порядка.

При получении порядка  $-pz$ , переполняющего разрядную сетку, формируются нулевой результат и признак исчезновения порядка.



**Пример.**

**Сложить  $4,63 + 4,63 = 9,26$**

$$\begin{array}{r} 0.100101000010000\dots\dots 0 \\ + \underline{0.100101000010000\dots\dots 0} \\ \color{red}{1.001010000100000\dots\dots 0} \end{array} \quad R_{cm}=127+2=129$$

**Возникло переполнение:**

**-необходимо произвести сдвиг вправо на 1 разряд,**

**-порядок увеличить на 1.**

$$R_{cm}=130, \quad q = 0.100101000010000\dots\dots 0$$

**Переведем число в десятичную СС:**

$$1001,0100001000 = -(1*2^0 + 1*2^3 + 1*2^{-2} + 1*2^{-7}) = 9,257812$$



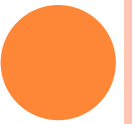
**Умножение** чисел с плавающей точкой выполняется в соответствии с формулой

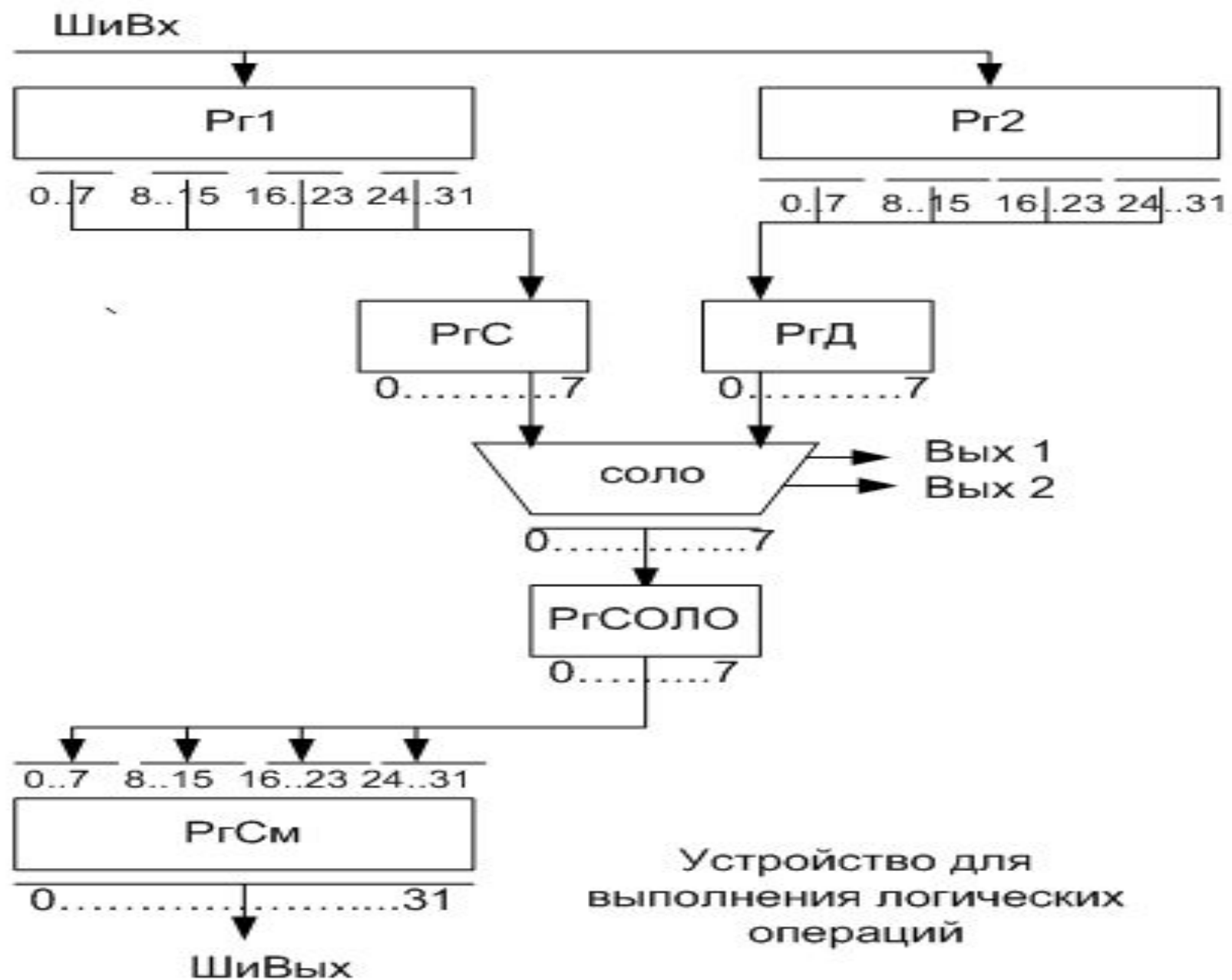
$$Z=X*Y=S^{P_x} q_x * S^{P_y} q_y = S^{(P_x + P_y)} q_x * q_y = S^{P_z} q_z$$

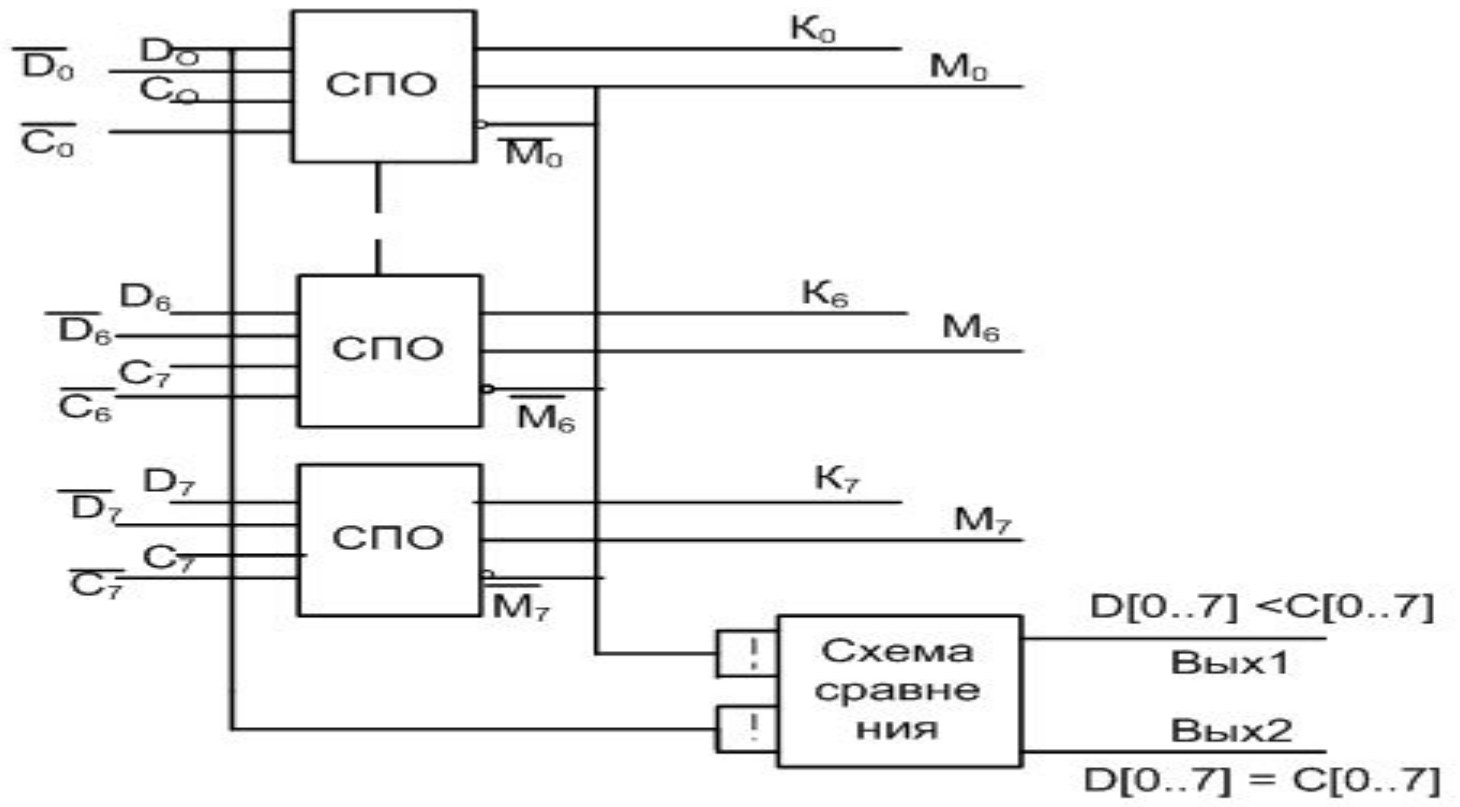
**Деление** чисел с плавающей точкой выполняется в соответствии с формулой

$$Z=X/Y=S^{P_x} q_x / S^{P_y} q_y = S^{(P_x - P_y)} q_x / q_y = S^{P_z} q_z$$









Комбинационная схема поразрядной обработки



**Вых 1**

**Вых 2**

**Результат  
сравнения**

**1**

**1**

**$D < C$**

**0**

**1**

**$D > C$**

**0**

**0**

**$D = C$**





На выходе  $K_i$  образуется конъюнкция

$$K_i = \overline{\overline{C_i} \vee \overline{D_i}} = C_i D_i$$

На выходе  $M_i$  формируется значение суммы по модулю 2

$$M_i = \overline{\overline{C_i} \overline{D_i} \vee C_i D_i} = C_i \oplus D_i$$

Выходы  $K_i$  и  $M_i$  соединены со входами РГСОЛО

$\text{РГСОЛО} := \text{РГС} \wedge \text{РГD}$  В РГСОЛО заносится значение состояния выходов  $K_i$

$\text{РГСОЛО} := \text{РГС} \oplus \text{РГD}$  В РГСОЛО заносится значение состояния выходов  $M_i$

Если выполняются 2 микрооперации, то в РГСОЛО заносится результат поразрядной операции ИЛИ

$$\text{РГСОЛО} := \text{РГС} \vee \text{РГD}$$



## Задание 1.

Преобразовать следующие числа с плавающей точкой одинарной точности из шестнадцатеричной в десятичную систему счисления.

Таблица 14. Исходные данные

№ вар	Исх.данные	№вар.	Исх.данные
1	4175C28FH	9	41A2041

## Задание 2.

Преобразовать следующие числа в формат стандарта IEEE с одинарной точностью. Результаты представить в восьми шестнадцатеричных разрядах.

Таблица 15. Исходные данные

№ варианта	Исх. данные	№ варианта	Исх. данные	№ варианта	Исх. данные
1	9	9	-12.45	17	-9.12
2	5.3210	6.2518	2.86		



*Задание 3.*

Сложить два числа с плавающей точкой на сумматоре прямого кода. Нормализовать результат и выразить его в шестнадцатеричной системе счисления.

Таблица 16. Исходные данные

№вар.	Исходные данные	№вар.	Исходные данные
1	3EE00000H+18800000H	13	FCFF0000H+ F79A8000H

*Задание 4.*

Выполнить умножение  $A*B$  и деление  $C/D$ , сложить полученные результаты, записать их в коротком формате.

Таблица 16. Исходные данные

№ вар	A	B	C	D	№ вар	A	B	C	D
1	15.360	0.38	1.78	132.7	13	201	1.220	9.9	0.49
2	0.57	195.3	0.67	0.89	14	56.37	42.81	0.9	3.57



№ вар.	операция	Условия нормализации мантиссы , особенности структуры, код
1	Сложение	$1 \geq  q_x  > 2$ , прямой код, схема СОЛО
2	Вычитание	$1 \geq  qx  > 2$ , прямой код, схема СОЛО
3	Умножение	$1 \geq  q_x  > 2$ , прямой код со старших разрядов
4	Деление	$1 \geq  q_x  > 2$ , доп. код без восстановления остатка.
5	Сложение	$1 \geq  qx  > 2$ , прямой код , использовать $\Sigma$ для выравнивания порядков
6	Вычитание	$1/s >  q_x  \geq 1$ , прямой код , использовать $\Sigma$ для выравнивания порядков
7	Деление	$1 \geq  q_x  > 2$ , прямой код, ускоренный алгоритм
8	Сложение	$1 \geq  qx  > 2$ , доп. код, схема СОЛО
9	Деление	$1 \geq  q_x  > 2$ , прямой код с восстановл. остатка
		$1 \geq  q_x  > 2$ , доп. код, схема СОЛО
10	Сложение	$1 \geq  q_x  > 2$ , прямой код, схема СОЛО
11	Вычитание	$1/s >  q_x  \geq 1$ , доп. код , использовать $\Sigma$ для выравнивания порядков
12	Умножение	$1/s >  q_x  \geq 1$ ., доп.код
13	Вычитание	$1/s >  q_x  \geq 1$ , прямой код , использовать $\Sigma$ для выравнивания порядков
14	Умножение	$1 \geq  qx  > 2$ , прямой код, М.С со старших разрядов
15	Деление	$1/s >  q_x  \geq 1$ ., доп.код
16	Умножение	$1 \geq  qx  > 2$ , прямой код, М.С со старших разрядов
17	Сложение	$1 \geq  q_x  > 2$ , прямой код, использовать $\Sigma$ для выравнивания порядков
18	Сложение	$1/s >  q_x  \geq 1$ , доп.код, использовать схему СОЛО