

Операции с плавающей точкой



ШКОЛА СИНТЕЗА
ЦИФРОВЫХ СХЕМ

Занятие №9
19 ноября 2022

ПРИ ПАРТНЕРСТВЕ



• • • •
• • • •
• • • •
YADRO • MP



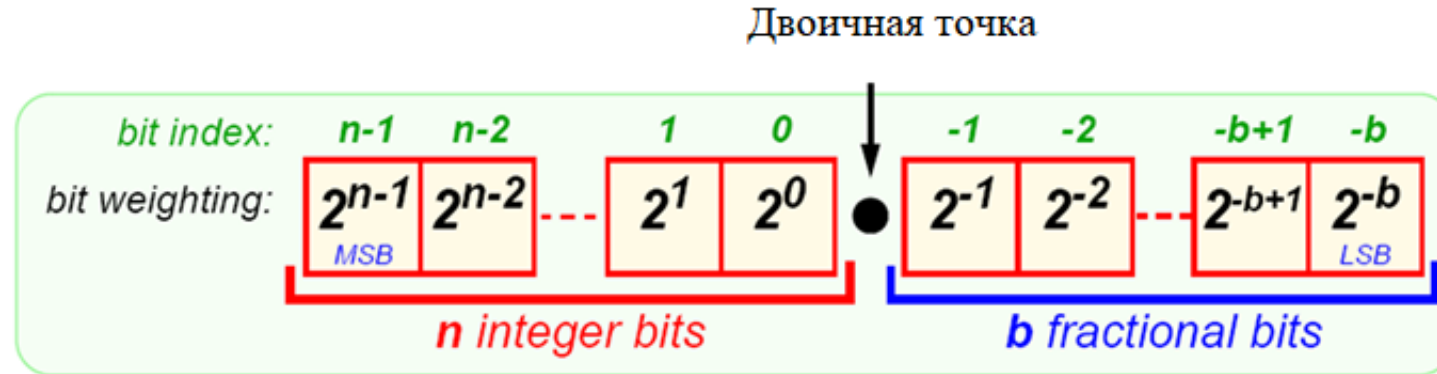


Илья Кудрявцев

Преподаватель Самарского университета

Дробные числа в формате с фиксированной точкой

Формат с фиксированной точкой предусматривает n целых бит и b дробных бит.



Вес дробных битов:

- $2^{-1} = 0.5$
- $2^{-2} = 0.25$
- $2^{-3} = 0.125....$

Дробные числа в формате с фиксированной точкой

Старший бит	ВЕС БИТА														Младший бит	
	-2^0	2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}	2^{-7}	2^{-8}	2^{-9}	2^{-10}	2^{-11}	2^{-12}	2^{-13}	2^{-14}	2^{-15}
Шестнадцатиричное представление	↓				↓											
7FFF	0111	1111	1111	1111									+0.999969			
0001	0000	0000	0000	0001									+0.000031			
0000	0000	0000	0000	0000									+0.000000			
FFFF	1111	1111	1111	1111									-0.000031			
8000	1000	0000	0000	0000									-1.000000			

Особенности формата с плавающей точкой

Проблемы формата с фиксированной точкой:

- Ограниченный динамический диапазон;
- Потеря точности (значащих цифр);

Достоинства формата с фиксированной точкой:

- Простота реализации;
- Использование тех же ресурсов, что и в целочисленных операциях

IEEE-754 : Формат с плавающей точкой (одинарная точность)

$$\text{NUMBER}_{10} = (-1)^S \times \underbrace{1.M}_{\text{Предполагается}} \times 2^{(E-127)} \quad \text{Смещение}$$



0	00000111	1100...00	
+	7	0.75	

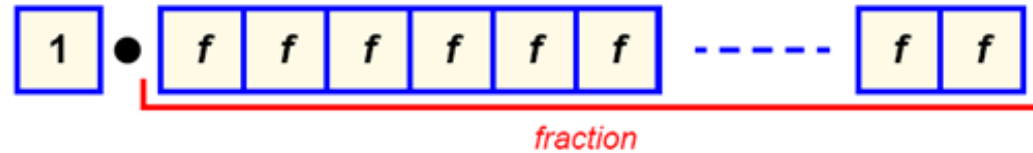
$$+ \underbrace{1.75}_{\text{Предполагается}} \times 2^{(7-127)} \quad \text{Смещение} = + 1,316554 \times 10^{-36}$$

1	10000001	0110...00	
-	129	0.375	

$$- \underbrace{1,375}_{\text{Предполагается}} \times 2^{(129-127)} \quad \text{Смещение} = - 5,500000$$

IEEE-754 : Формат с плавающей точкой

Согласно IEEE 754, мантисса должна быть нормализованной, т.е. представленной в следующей форме (целочисленный бит должен всегда иметь значение 1)



Целая часть мантиссы указывается неявно, т.е. представляет собой «скрытый бит»

Мантисса умножается на знаковый бит, т.е. на +1 или -1

$s = 0:$ $(-1)^s = (-1)^0 = 1$ Положительные числа

$s = 1:$ $(-1)^s = (-1)^1 = -1$ Отрицательные числа

IEEE-754 : Формат с плавающей точкой

IEEE 754 определяет четыре варианта формата / точности.

Чаще всего используется формат одинарной точности (**Single**) и формат двойной точности (**Double**).

Length (bits)	Format			
	Single	Single Extended	Double	Double Extended
Significand	24	≥ 32	53	≥ 64
Exponent	8	≥ 11	11	≥ 15
Total	32	≥ 43	64	≥ 79

В форматах одинарной и двойной точности используются, соответственно, 23-х разрядная и 52-разрядная мантисса

Для формата с одинарной точностью требуется:

$$23 \text{ бита мантиссы} + 8 \text{ бит порядка} + 1 \text{ бит знака} = 32 \text{ бит.}$$

Для формата с двойной точностью требуется:

$$52 \text{ бита мантиссы} + 11 \text{ бит порядка} + 1 \text{ бит знака} = 64 \text{ бит.}$$

IEEE-754 : Формат с плавающей точкой

- 16-разрядная с фиксированной точкой:
 - ♦ $2^{16} = 65536$ возможных значений
- 32-разрядная с плавающей точкой:
 - ♦ Наибольшее значение: $\pm 6.8 \times 10^{38}$, в стандарте IEEE-754: $\pm 3.4 \times 10^{38}$
 - ♦ Наименьшее значение: $\pm 5.9 \times 10^{-39}$, в стандарте IEEE-754: $\pm 1.2 \times 10^{-38}$
- Расширенная точность (40-бит: знак + 8-битная экспонента + 31-битная мантисса)
- Двойная точность: (64-бит.: знак + 11-бит. эксп.+ 52-бит. мантисса)
- 32-разрядная с плавающей точкой:
 - ♦ Более высокая точность
 - ♦ Большой динамический диапазон
 - ♦ Проще в программировании

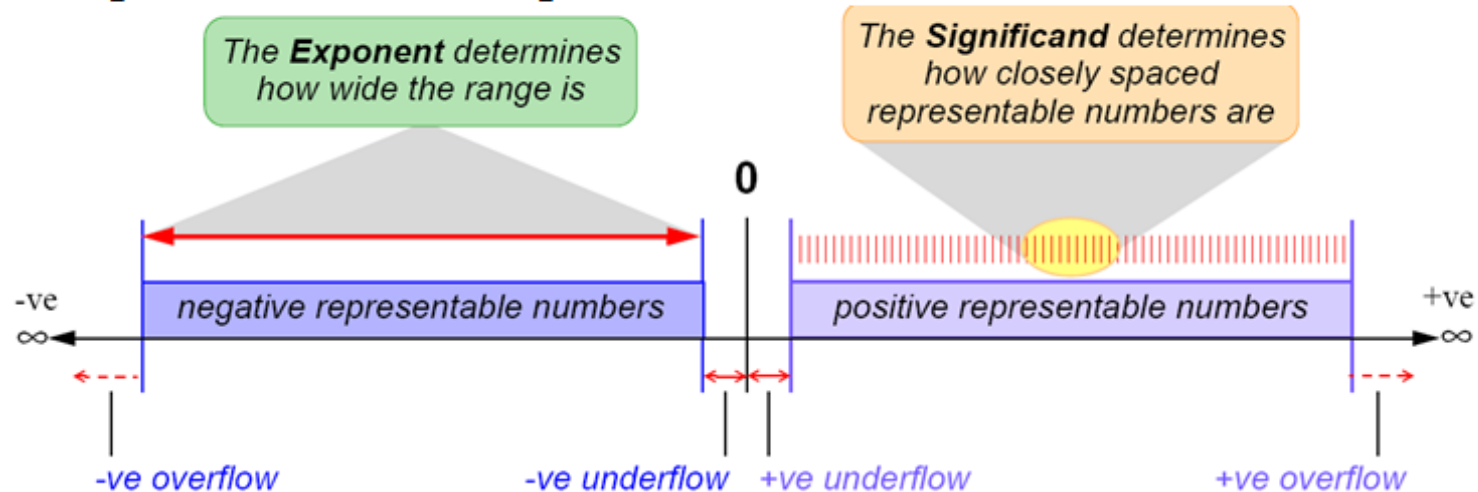
IEEE-754 : Формат с плавающей точкой

Согласно IEEE-754 должны поддерживаться следующие операции с плавающей точкой:

- Сложение, вычитание и умножение
- Деление, вычисление квадратного корня и остатка
- Сравнения
- Операции округления
- Преобразования между форматами
- Преобразования между форматами с плавающей точкой и целочисленными

IEEE-754 : Формат с плавающей точкой

- Мантисса и порядок определяют различные свойства чисел в формате с плавающей точкой
- Мантисса определяет разрешение
- Порядок определяет динамический диапазон числа
- Аналогично арифметике с фиксированной точкой, возможны переполнение и потеря точности



IEEE-754 : Формат с плавающей точкой (32 бит)

<i>Table 1–2. Representation of Special-Case Numbers in IEEE 754-1985 Standard</i>			
Meaning	Sign Field	Exponent Field	Mantissa Field
Zero	0/1	All 0's	All 0's
Positive Denormalized	0	All 0's	Non-zero
Negative Denormalized	1	All 0's	Non-zero
Positive Infinity	0	All 1's	All 0's
Negative Infinity	1	All 1's	All 0's
Not a Number (NaN)	0/1	All 1's	Non-zero

NaN	Zero	Infinity	Denormalized
<ul style="list-style-type: none"> • $(+\infty) + (-\infty) = \text{NaN}$ • $\text{NaN} + x = \text{NaN}$ 	<ul style="list-style-type: none"> • $(+a) + (-a) = +0$ • $(+0) + (-0) = +0$ • $(-0) + (-0) = -0$ 	<ul style="list-style-type: none"> • <i>overflow</i>: $(a) + (a) = \text{inf}$ • $(\text{inf}) + x = \text{inf}$ • $(\text{inf}) + (\text{inf}) = \text{inf}$ 	<ul style="list-style-type: none"> • <i>underflow</i>: $(a) + (-b) = \text{denorm}$

IEEE-754 : Формат с плавающей точкой (32 бит)

Tools & Thoughts

IEEE-754 Floating Point Converter

Translations: [de](#)

This page allows you to convert between the decimal representation of numbers (like "1.02") and the binary format used by all modern CPUs (IEEE 754 floating point).

[illegible][illegible]

IEEE-754 : Формат с плавающей точкой (32 бит)

[illegible]

IEEE 754 Converter (JavaScript), V0.22

		Sign	Exponent	Mantissa
Value:		+1	2^{128}	1.0000001192092896
Encoded as:		0	255	1
Binary:		<input type="checkbox"/>	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/>
Decimal representation		<input type="text" value="nan"/>		
Value actually stored in float:		<input type="text" value="NaN"/>		
Error due to conversion:		<input type="text"/>		
Binary Representation		<input type="text" value="01111111100000000000000000000001"/>		
Hexadecimal Representation		<input type="text" value="0x7f800001"/>		

+1

-1

IEEE-754 : Формат с плавающей точкой (32 бит)

IEEE 754 Converter (JavaScript), V0.22

	Sign	Exponent	Mantissa
Value:	+1	2 ⁻¹²⁶ (denormalized)	1.1920928955078125e-07 (denormalized)
Encoded as:	0	0	1
Binary:	<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/>
Decimal representation	<input type="text" value="1.40129846432e-45"/>		
Value actually stored in float:	<input type="text" value="1.401298464324817070923729583289916131280261941876515771751"/>		
Error due to conversion:	<input type="text"/>		
Binary Representation	<input type="text" value="00000000000000000000000000000001"/>		
Hexadecimal Representation	<input type="text" value="0x00000001"/>		

+1
-1

IEEE 754 Converter (JavaScript), V0.22

	Sign	Exponent	Mantissa
Value:	-1	2^{-126} (denormalized)	0.0 (denormalized)
Encoded as:	1	0	0
Binary:	<input checked="" type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
Decimal representation	<input type="text" value="-0.0"/>		
Value actually stored in float:	<input type="text" value="-0"/>		
Error due to conversion:	<input type="text"/>		
Binary Representation	<input type="text" value="10000000000000000000000000000000"/>		
Hexadecimal Representation	<input type="text" value="0x80000000"/>		

+1 -1

IEEE-754 : Формат с плавающей точкой (32 бит)

IEEE 754 Converter (JavaScript), V0.22

	Sign	Exponent	Mantissa
Value:	+1	2^{-4}	1.6000000023841858
Encoded as:	0	123	5033165
Binary:	<input type="checkbox"/>	<input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/>
You entered	<input type="text" value="0.1"/>		
Value actually stored in float:	<input type="text" value="0.100000001490116119384765625"/>		
Error due to conversion:	<input type="text" value="1.490116119384765625E-9"/>		
Binary Representation	<input type="text" value="0011110111001100110011001101"/>		
Hexadecimal Representation	<input type="text" value="0x3dcccccd"/>		

☒ +1 ☒ -1

IEEE 754 Converter (JavaScript), V0.22

	Sign	Exponent	Mantissa
Value:	+1	2^{127}	1.7632415294647217
Encoded as:	0	254	6402534
Binary:	<input type="checkbox"/>	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/>	<input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input checked="" type="checkbox"/> <input checked="" type="checkbox"/> <input type="checkbox"/>
You entered	<input type="text" value="3.0E+38"/>		
Value actually stored in float:	<input type="text" value="300000000549775575777803994281145270272"/>		
Error due to conversion:	<input type="text" value="5.497755757778039942811452703E+29"/>		
Binary Representation	<input type="text" value="0111111011000011011000111100110"/>		
Hexadecimal Representation	<input type="text" value="0x7f61b1e6"/>		

☒ +1 ☒ -1

IEEE-754 : Исключения в операциях с плавающей точкой

Исключение	Примечания
Переполнение	Результат $\pm\infty$
Потеря точности	Результат равен нулю или денормализован
Деление на нуль	Результат $\pm\infty$
Некорректная операция	Результат NaN (не число)
Неточность	Требуется округление

IEEE-754 : Операции с некорректным результатом

Операция	Примечания
Сложение/вычитание	Операции типа $\infty \pm \infty$
Умножение	Операции типа $0 \times \infty$
Деление	Операции типа $0/0$ или ∞/∞
Вычисление остатка	Операции типа $x \text{ REM } 0$ или $\infty \text{ REM } y$
Вычисление квадратного корня	Квадратный корень из отрицательного числа

Денормализованные числа

Побитовое представление: **s_00000000_xxxxxxxxxxxxxxxxxxxxxxxxxx** (мантиса != 0)

Представление числа в
десятичном виде:

$$F = (-1)^s 2^{-126} (M / 2^{23})$$

Образуются при вычитании малых чисел, пример:

$$a = 0_00000100_000000000000000000000000 \sim 2^{-123}$$

$$b = 1_00000011_11010000100111000001111 \sim -1.81488 * 2^{-124}$$

$$1) \quad a = 1.000000000000000000000000 * 2^{-123}$$

$$b = -0.11101000010011100000111|1 * 2^{-123}$$

$$2) \quad \text{sum} = 0.00010111101100011111000|1 * 2^{-123} = 0.10111101100011111000100 * 2^{-126}$$

$$\text{sum} = 0_00000000_10111101100011111000100 \sim 0.74047 * 2^{-126}$$

Усечение и округление

- **Усечение** и **округление** представляют собой операции уменьшения точности представления числа для сокращения размера мантиссы
- Рассмотрим пример в десятичной арифметике
- Предположим, что нужно сократить 27.3147 до трех цифр после десятичной точки
- После усечения получим 27.314, а после округления 27.315, так как операции отличаются

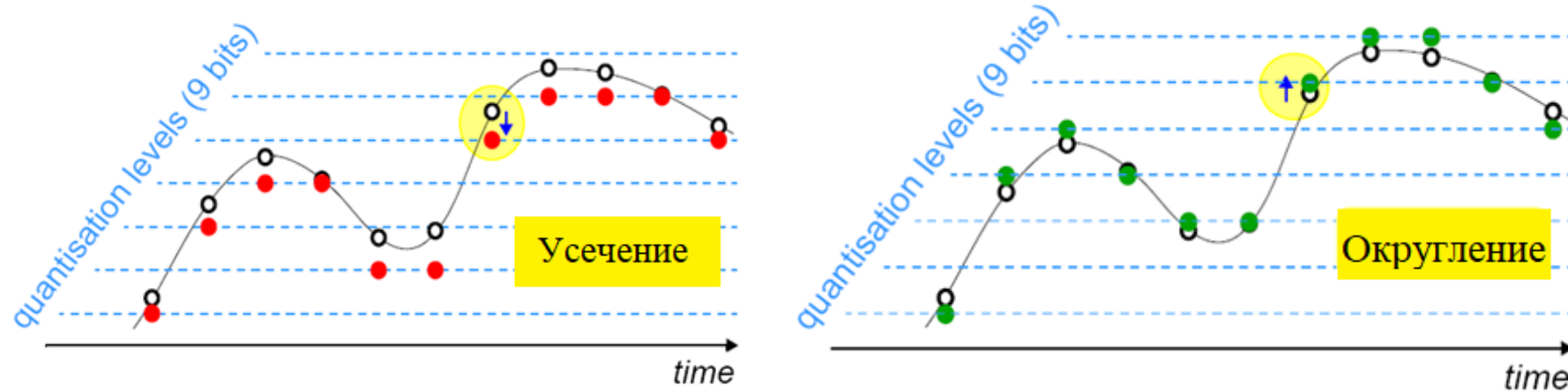
$$\begin{array}{r} 24.3147 \\ \downarrow \\ 24.314 \end{array}$$

Усечение означает отбрасывание младших цифр мантиссы

$$\begin{array}{r} 24.3147 \\ + 00.0005 \\ \hline 24.3152 \\ \downarrow \\ 24.315 \end{array}$$

Округление означает добавление половины младшей значащей цифры перед отбрасыванием

Усечение и округление



Операция округления обеспечивает более высокую точность, чем операция усечения, но требует дополнительной операции сложения.

Округление

Мантисса после сдвига: 1100100110011001010101010101101111

L - last bit - последний бит, который помещается в размер мантиссы.

G - guard bit - первый не помещающийся бит

R - round bit - второй не помещающийся бит

S - sticky bit - логическая сумма всех оставшихся битов

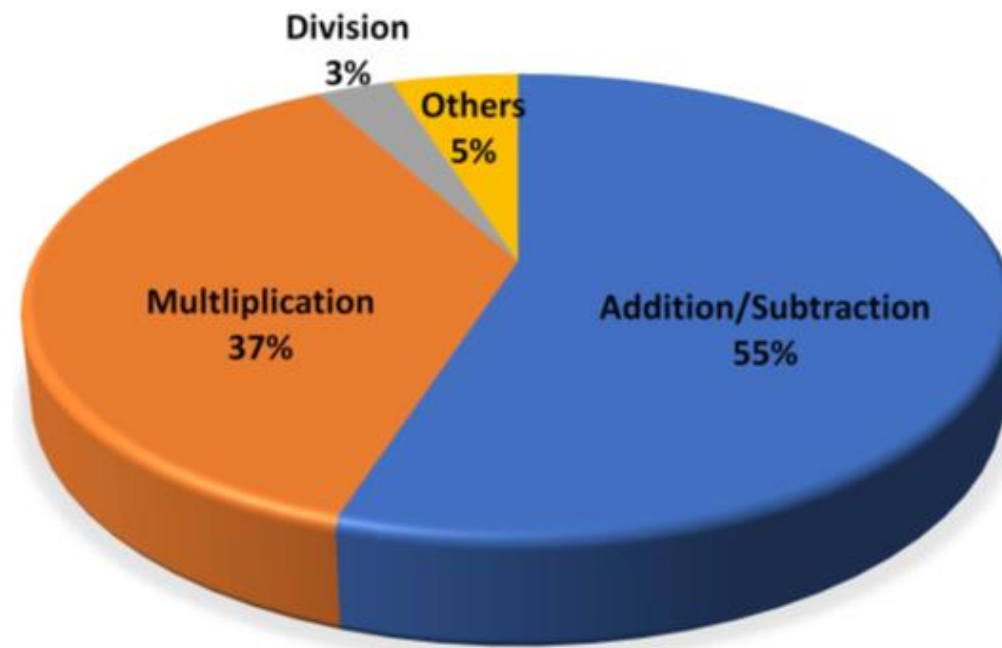
Округление происходит в случае, если условие **G & (L | R | S)** верно.

Таблица 4.1.

L	G	R	St	C
X	0	X	X	0
1	1	X	X	1
X	1	1	X	1
X	1	X	1	1

$$C_r = G(R + St + L) = 1,$$

Типичные операции IEEE-754



Asynchronous Floating-Point Adders and Communication Protocols: A Survey
Pallavi Srivastava 1, Edwin Chung 1 and Stepan Ozana 2

Алгоритм операции умножения

- Шаг 1 : Умножение мантисс :

$$SIG_R = SIG_A \times SIG_B = 1.375 \times 1.5 = 2.0625$$

- Шаг 2 : Сложение порядков

$$E_R = E_A + E_B = 5 + 7 = 12$$

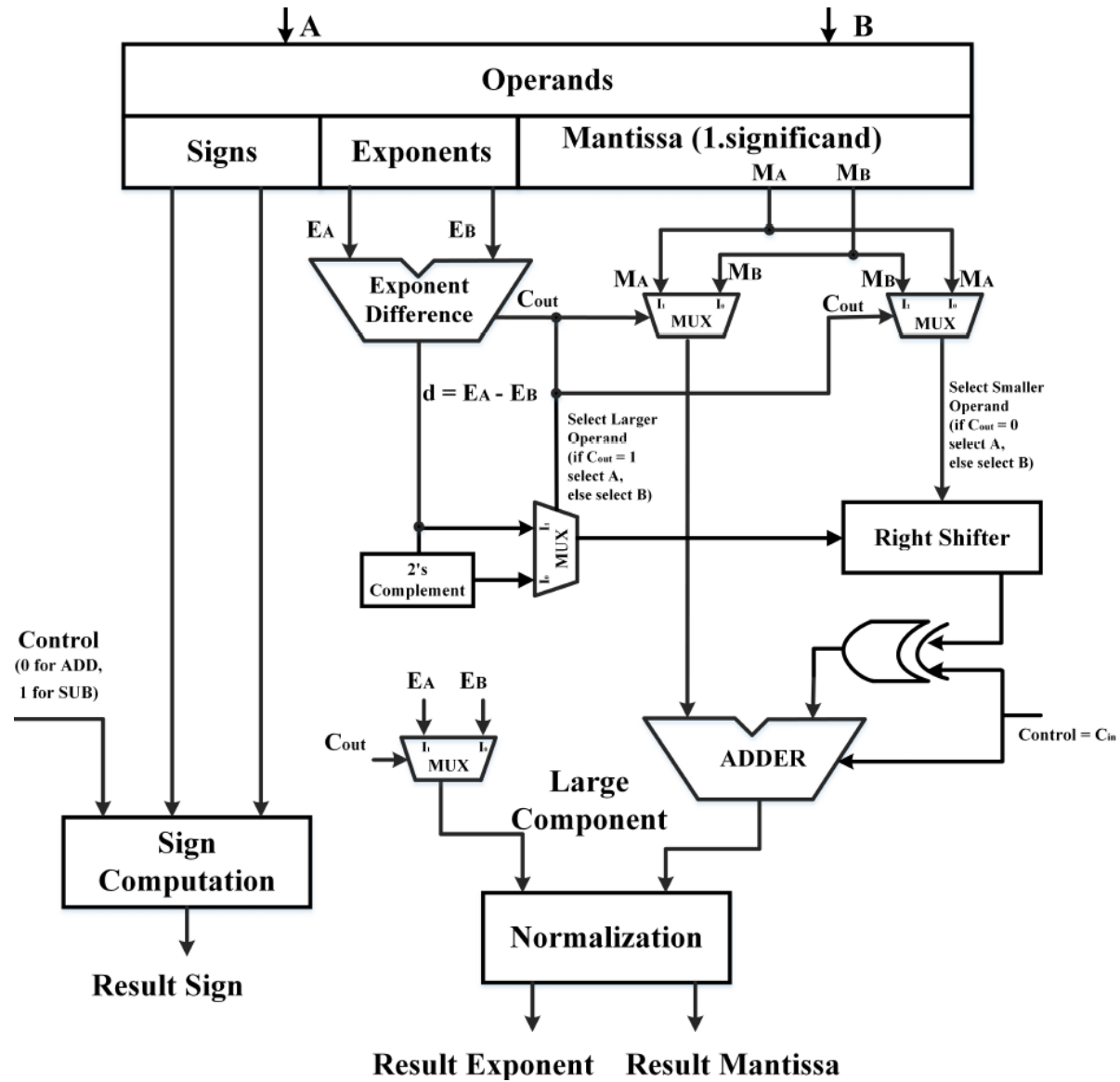
- Шаг 3 : Нормализация результата путем коррекции мантиссы и порядка.
(это необходимо, так как мантисса $SIG_R > 2$)

$$SIG_{R_{norm}} = \frac{SIG_R}{2} = \frac{2.0625}{2} = 1.03125$$

$$E_{R_{norm}} = E_R + 1 = 12 + 1 = 13$$

- Шаг 4 : Определение порядка результата $S_R = S_A \text{ xor } S_B = 1$

Операция сложения

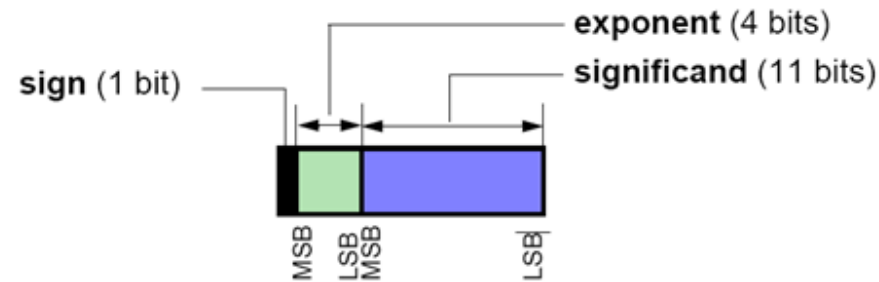


Блочная плавающая точка

- Блочная плавающая точка – это технология работы с блоком чисел с одинаковым порядком (но разными мантиссами);
- Операции осуществляются над элементами блока без необходимости постоянной нормализации;
- Порядок корректируется только в случае сильного увеличения или уменьшения чисел;
- Величина порядка определяется числом максимальной величины;
- За счет упрощения техники применение блочной техники требует значительно меньшего объема ресурсов, чем стандартные операции с плавающей точкой;
- Технология блочной плавающей точки хорошо подходит к алгоритмам типа БПФ и реализуется в сигнальных процессорах.

Нестандартные форматы

- Не существует жесткого требования обязательного соответствия IEEE-754
- Анализируя требования приложения к точности и динамическому диапазону, можно подобрать соответствующие параметры формата с плавающей точкой
- Пример нестандартного формата с 4-х разрядным порядком и 11-разрядной мантиссой (16 разрядов с учетом знака)



- По сравнению со стандартным форматом обеспечивается более широкий динамический диапазон при незначительном снижении точности

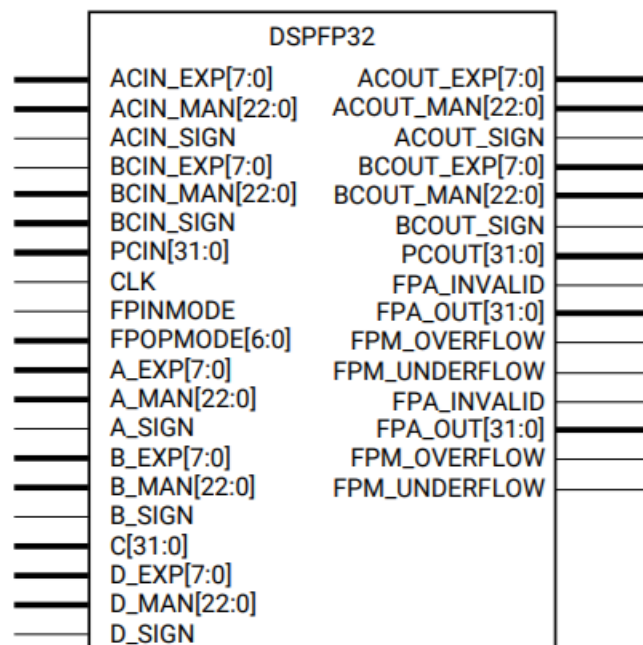
Поддержка операций с плавающей точкой в Versal

DSPFP32

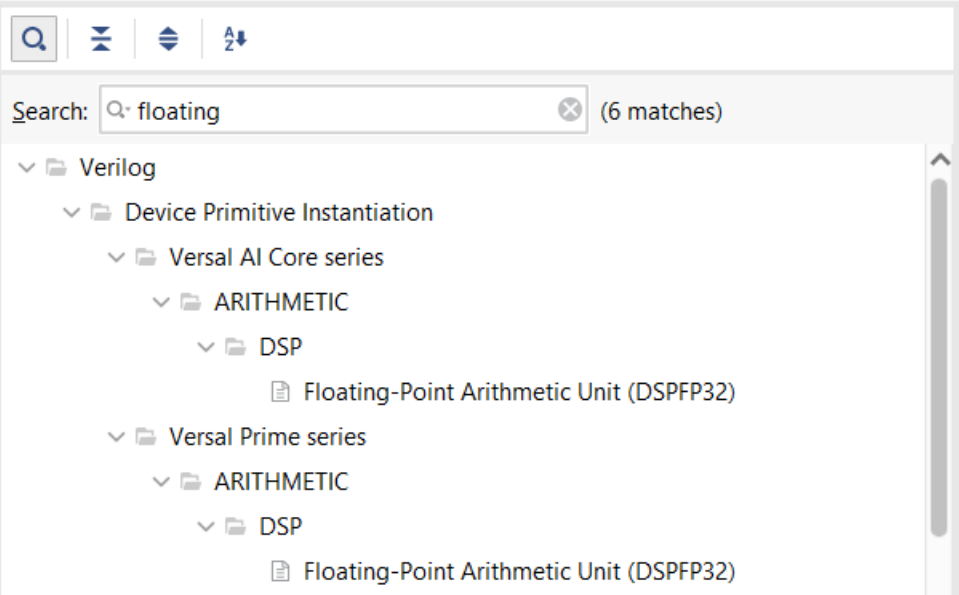
Primitive: The DSPFP32 consists of a floating-point multiplier and a floating-point adder with separate outputs.

PRIMITIVE_GROUP: **ARITHMETIC**

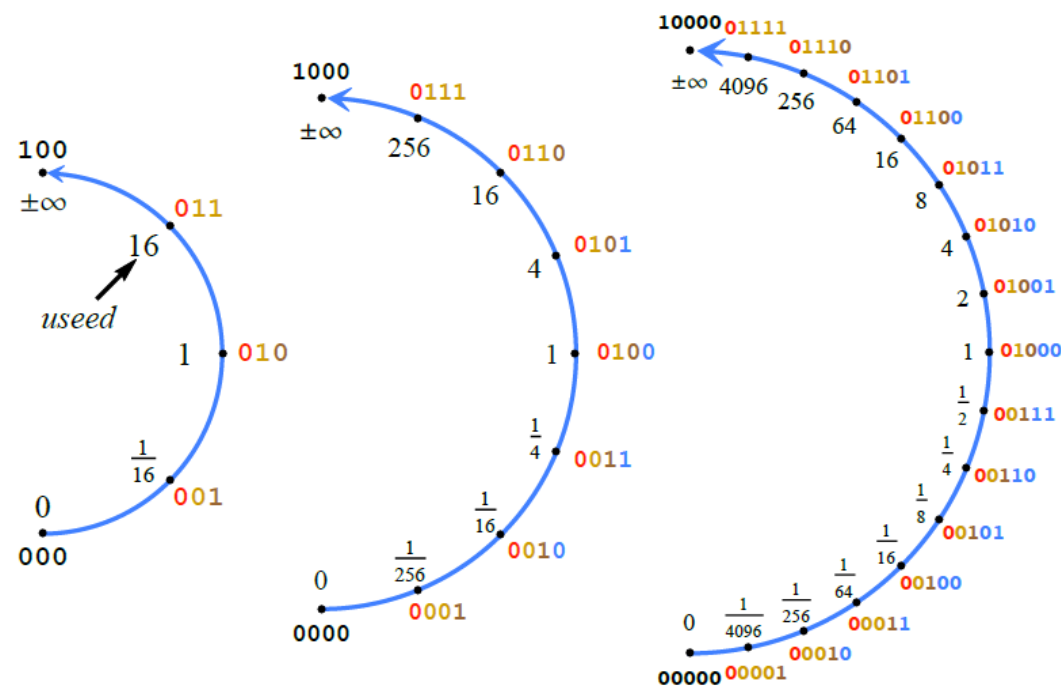
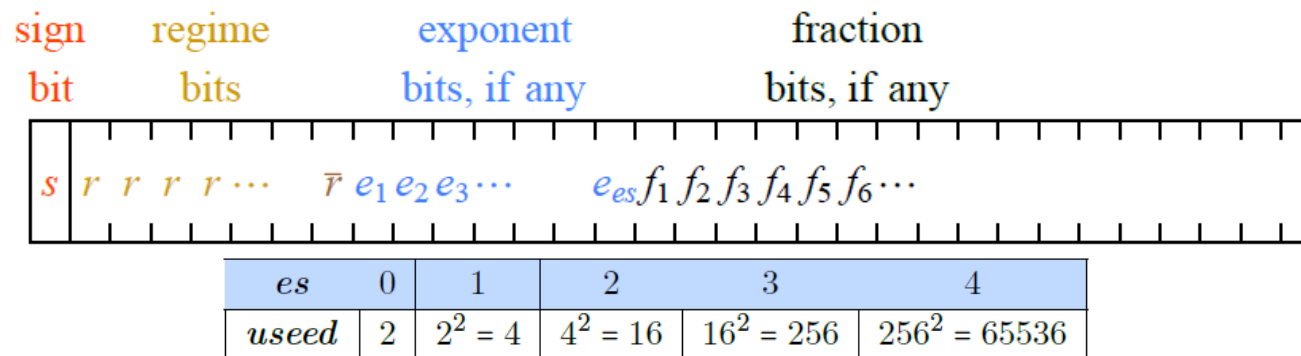
PRIMITIVE_SUBGROUP: DSP



Templates



Posits



Beating Floating Point at its Own Game: Posit Arithmetic
John L. Gustafson , Isaac Yonemoto

Posits

$$x = \begin{cases} 0, & p = 0, \\ \pm\infty, & p = -2^{n-1}, \\ \text{sign}(p) \times \text{useed}^k \times 2^e \times f, & \text{all other } p. \end{cases}$$

