# Double Precision Floating Pt.

## 64 bits   (IEEE 754)

52 bits

11 bits   52 51

63 62

0

Mantissa Bits

sign bit

Exponent Bits

$$\text{real value (base 10)} = (-1)^{\text{sign}} \left(1 + \sum_{i=1}^{52} b_{52-i} \, 2^{-i}\right) \times 2^{e-1023}$$

## Largest #

```
0 | 11111111110 | 1111 - - - - - - - - - (1111)
```

$2^{11}-1-1$

$\approx 1$

$$\angle \quad = \quad \ldots$$

$$= 2046$$

$$r_{max} = 2 \times 2^{2046-1023}$$

$$= 2 \times 2^{1023}$$

$$= 2^{1024}$$

$$\boxed{r_{max} = 1.798 \times 10^{308}}$$

## Smallest #

$$\boxed{0} \; \underbrace{0000\;0000\;001}_{=\;1} \; \underbrace{0 \;-\;-\;-\;-\;-\;-\;-\;-\;0}_{=\;1}$$

$$r_{min} = 1.0 \times 2^{1-1023}$$

$$= 1.0 \times 2^{-1022}$$

$$\boxed{r_{min} = 2.23 \times 10^{-308}}$$

Smallest difference between 2 Numbers.

$$\boxed{.000 \cdots\cdots\cdots\cdots\cdots\cdots .1}$$

$$= 2^{-52} = \boxed{2.22 \times 10^{-16}}$$

This is why we say that double precision is accurate to 15 decimal places.

---

How is $\pi$ stored in double precision?

$$\pi \approx 3.141592653589793$$

(to 16 digits)



$$r = (-1)^s \left(1 + \sum_{i=1}^{52} b_{52-i} \, 2^{-i}\right) \times 2^{e-1023}$$

$e - 1023 = 1 \quad \therefore e = 1024$

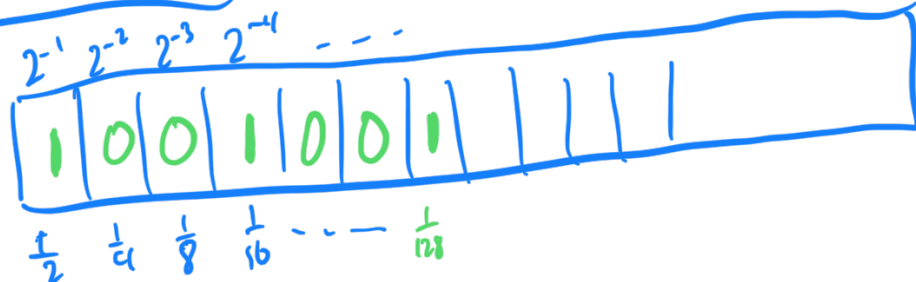$$= \boxed{\phantom{xxxx}} \times 2^{\boxed{1}}$$

$\mathbb{1} = \underbrace{\qquad}$     exponent ~~~~

$(1 \leftrightarrow 2)$

$$\pi = \boxed{\phantom{xxx}} \times 2^1$$

$$= \frac{\pi}{2} = 1 + 0.570796327\ldots\ldots$$

bits of mantissa

$2^{-1}$   $2^{-2}$   $2^{-3}$   $2^{-4}$   — — —

| 1 | 0 | 0 | 1 | 0 | 0 | 1 | | | | | |

$\frac{1}{2}$   $\frac{1}{4}$   $\frac{1}{8}$   $\frac{1}{16}$   — —   $\frac{1}{128}$

$0.5707496327$

$\sim 0.5$ _____

$-\frac{1}{16}$ _____

$0.008296327$

$-\frac{1}{128}$ _____      etz.

$0.000483827\ldots$

| 1 | 10000 000000 | 1 001 001 0000 1111111   — — — —   11000 |

## Double-precision examples  [ edit ]

0 01111111111 0000000000000000000000000000000000000000000000000000$_2$ ≙ 3FF0 0000 0000 0000$_{16}$ ≙ $+2^0 \times 1 = 1$

0 01111111111 0000000000000000000000000000000000000000000000000001$_2$ ≙ 3FF0 0000 0000 0001$_{16}$ ≙ $+2^0 \times (1 + 2^{-52}) \approx 1.0000000000000002$, the smallest number > 1

0 01111111111 0000000000000000000000000000000000000000000000000010$_2$ ≙ 3FF0 0000 0000 0002$_{16}$ ≙ $+2^0 \times (1 + 2^{-51}) \approx 1.0000000000000004$

0 10000000000 0000000000000000000000000000000000000000000000000000$_2$ ≙ 4000 0000 0000 0000$_{16}$ ≙ $+2^1 \times 1 = 2$

1 10000000000 0000000000000000000000000000000000000000000000000000$_2$ ≙ C000 0000 0000 0000$_{16}$ ≙ $-2^1 \times 1 = -2$

0 10000000000 1000000000000000000000000000000000000000000000000000$_2$ ≙ 4008 0000 0000 0000$_{16}$ ≙ $+2^1 \times 1.1_2 = 11_2 = 3$

0 10000000001 0000000000000000000000000000000000000000000000000000$_2$ ≙ 4010 0000 0000 0000$_{16}$ ≙ $+2^2 \times 1 = 100_2 = 4$

0 10000000001 0100000000000000000000000000000000000000000000000000$_2$ ≙ 4014 0000 0000 0000$_{16}$ ≙ $+2^2 \times 1.01_2 = 101_2 = 5$

0 10000000001 1000000000000000000000000000000000000000000000000000$_2$ ≙ 4018 0000 0000 0000$_{16}$ ≙ $+2^2 \times 1.1_2 = 110_2 = 6$

0 10000000011 0111000000000000000000000000000000000000000000000000$_2$ ≙ 4037 0000 0000 0000$_{16}$ ≙ $+2^4 \times 1.0111_2 = 10111_2 = 23$

0 01111111000 1000000000000000000000000000000000000000000000000000$_2$ ≙ 3F88 0000 0000 0000$_{16}$ ≙ $+2^{-7} \times 1.1_2 = 0.00000011_2 = 0.01171875$ (3/256)

0 00000000000 0000000000000000000000000000000000000000000000000001$_2$ ≙ 0000 0000 0000 0001$_{16}$ ≙ $+2^{-1022} \times 2^{-52} = 2^{-1074} \approx 4.9406564584124654 \times 10^{-324}$ (Min. subnormal positive double)

0 00000000000 1111111111111111111111111111111111111111111111111111$_2$ ≙ 000F FFFF FFFF FFFF$_{16}$ ≙ $+2^{-1022} \times (1 - 2^{-52}) \approx 2.2250738585072009 \times 10^{-308}$ (Max. subnormal double)

0 00000000001 0000000000000000000000000000000000000000000000000000$_2$ ≙ 0010 0000 0000 0000$_{16}$ ≙ $+2^{-1022} \times 1 \approx 2.2250738585072014 \times 10^{-308}$ (Min. normal positive double)

0 11111111110 1111111111111111111111111111111111111111111111111111$_2$ ≙ 7FEF FFFF FFFF FFFF$_{16}$ ≙ $+2^{1023} \times (1 + (1 - 2^{-52})) \approx 1.7976931348623157 \times 10^{308}$ (Max. Double)

0 00000000000 0000000000000000000000000000000000000000000000000000$_2$ ≙ 0000 0000 0000 0000$_{16}$ ≙ $+0$

1 00000000000 0000000000000000000000000000000000000000000000000000$_2$ ≙ 8000 0000 0000 0000$_{16}$ ≙ $-0$

0 11111111111 0000000000000000000000000000000000000000000000000000$_2$ ≙ 7FF0 0000 0000 0000$_{16}$ ≙ $+\infty$ (positive infinity)

1 11111111111 0000000000000000000000000000000000000000000000000000$_2$ ≙ FFF0 0000 0000 0000$_{16}$ ≙ $-\infty$ (negative infinity)

0 11111111111 0000000000000000000000000000000000000000000000000001$_2$ ≙ 7FF0 0000 0000 0001$_{16}$ ≙ NaN (sNaN on most processors, such as x86 and ARM)

0 11111111111 1000000000000000000000000000000000000000000000000001$_2$ ≙ 7FF8 0000 0000 0001$_{16}$ ≙ NaN (qNaN on most processors, such as x86 and ARM)

0 11111111111 1111111111111111111111111111111111111111111111111111$_2$ ≙ 7FFF FFFF FFFF FFFF$_{16}$ ≙ NaN (an alternative encoding of NaN)

0 01111111101 0101010101010101010101010101010101010101010101010101$_2$ = 3FD5 5555 5555 5555$_{16}$ ≙ $+2^{-2} \times (1 + 2^{-2} + 2^{-4} + \ldots + 2^{-52}) \approx \frac{1}{3}$

0 10000000000 1001001000011111101101010100010001000010110100011000$_2$ = 4009 21FB 5444 2D18$_{16}$ ≈ pi