# IEEE Rounding

$0, 1, 2, 3, 4$     round $\Downarrow$

$6, 7, 8, 9$     round $\Uparrow$

$5$     round to the nearest even digit

$1.5 \implies 2.0$

$2.5 \implies 2.0$

Remember, there are an infinite number of values between $0.0 \rightarrow 1.0$. We can only represent $2^{53}$ with double precision FP

example

$$2.56_{16} \times 10^4 + 2.34_{16} \times 10^2$$

We perfor

$$
\begin{array}{r}
2.3400 \\
+ \quad 0.0256 \\
\hline
2.3456
\end{array}
$$

So, $2.3656 \times 10^2$

We must now round to 2 significant digits to the right of the decimal.

$$00 \Rightarrow 49 \quad \text{round} \downarrow$$
$$51 \Rightarrow 99 \quad \text{round} \uparrow$$
$$50 \qquad \text{to nearest even}$$

In this case $2.37 \times 10^2$

# Misc Trivia ?

IEEE 754 defines 32, 64, & 128 bit FP formats

C# has a decimal type that is 128 bits

C/C++   float is 32 bits
         double is 64 bits

Check on 128 bits ?

1/10 in binary is

$$0.00011001100110011\ldots$$

1/10 is rational in decimal, but irrational in binary.

Most monetary calculations are performed using binary coded

performed using binary coded decimals.

When performing binary floating point operations, the hardware will (typically) provide two additional bits that are named <u>guard</u> and <u>round</u>

With 5 significant digits

$$0.1011 \times 2^0$$
$$+ \ 0.1111 \times 2^{-2}$$

<u>5 sig digits</u> ← guard
← round

$$0.1011$$
$$+ \quad 0.001111$$
$$\overline{\phantom{xxxxxxxxxxxxxxxxxxxxx}}$$
$$0.110\underline{11} \quad \text{round up since this is } 3_{16}$$
$$0.1111 \times 2^0$$

patterns
00 & 01 round ↓
& 10 & 11 round ↑