

## Using The Conversion Procedure

- Convert 2.625 to our 8-bit floating point format.
  - A. The integral part is easy,  $2_{10} = 10_2$ . For the fractional part:

$$\begin{array}{rcl}
 0.625 \times 2 = 1.25 & \boxed{1} & \text{Generate 1 and continue with the rest.} \\
 0.25 \times 2 = 0.5 & \boxed{0} & \text{Generate 0 and continue.} \\
 0.5 \times 2 = 1.0 & \boxed{1} & \text{Generate 1 and nothing remains.}
 \end{array}$$

- B. So  $0.625_{10} = 0.101_2$ , and  $2.625_{10} = 10.101_2$ .
- C. Add an exponent part:  $10.101_2 = 10.101_2 \times 2^0$ .
- D. Normalize:  $10.101_2 \times 2^0 = 1.0101_2 \times 2^1$ .
- E. Mantissa: 0101
- F. Exponent:  $1 + 3 = 4 = 100_2$ .
- G. Sign bit is 0.

The result is 0 100 0101.

- Convert -4.75 to our 8-bit floating point format.
  - a. The integral part is  $4_{10} = 100_2$ . The fractional:

$$\begin{array}{rcl}
 0.75 \times 2 = 1.5 & \boxed{1} & \text{Generate 1 and continue with the rest.} \\
 0.5 \times 2 = 1.0 & \boxed{1} & \text{Generate 1 and nothing remains.}
 \end{array}$$

- b. So  $4.75_{10} = 100.11_2$ .
- c. Normalize:  $100.11_2 = 1.0011_2 \times 2^2$ .
- d. Mantissa is 0011, exponent is  $2 + 3 = 5 = 101_2$ , sign bit is 1.

So -4.75 is = 1 101 0011

Convert 0.40625 to our 8-bit floating point format.

. Converting:

$$\begin{array}{rcl}
 0.40625 \times 2 = 0.8125 & \boxed{0} & \text{Generate 0 and continue.} \\
 0.8125 \times 2 = 1.625 & \boxed{1} & \text{Generate 1 and continue with the rest.} \\
 0.625 \times 2 = 1.25 & \boxed{1} & \text{Generate 1 and continue with the rest.} \\
 0.25 \times 2 = 0.5 & \boxed{0} & \text{Generate 0 and continue.} \\
 0.5 \times 2 = 1.0 & \boxed{1} & \text{Generate 1 and nothing remains.}
 \end{array}$$

- a. So  $0.40625_{10} = 0.01101_2$ .
- b. Normalize:  $0.01101_2 = 1.101_2 \times 2^{-2}$ .
- c. Mantissa is 1010, exponent is  $-2 + 3 = 1 = 001_2$ , sign bit is 0.

So 0.40625 is 0 001 1010