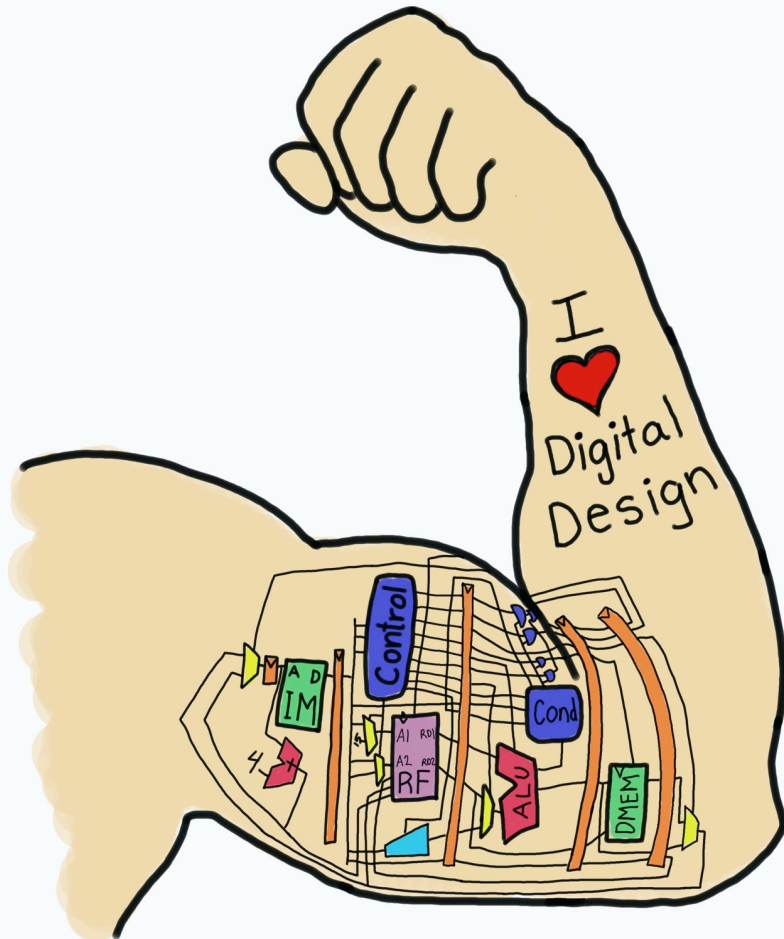


# E85 Digital Design & Computer Engineering



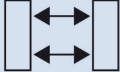
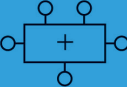

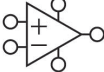




## Lecture 9: Shift, Mult, Div Fixed & Floating Point

**HARVEY  
MUDD  
COLLEGE**

# Lecture 9

- Shifters
- Multipliers
- Dividers
- Fixed Point Numbers
- Floating Point Numbers

Application Software	>"hello world!"
Operating Systems	
Architecture	
Micro-architecture	
Logic	
Digital Circuits	
Analog Circuits	
Devices	
Physics	



# Shifters

**Logical shifter:** shifts value to left or right and fills empty spaces with 0's

- Ex: **11001** >> 2 =
- Ex: **11001** << 2 =

**Arithmetic shifter:** same as logical shifter, but on right shift, fills empty spaces with the old most significant bit (msb)

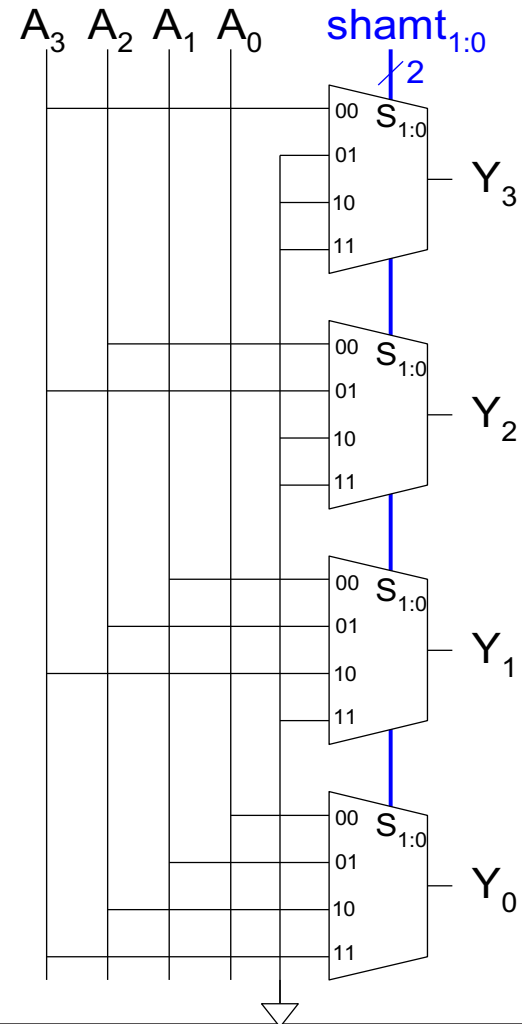
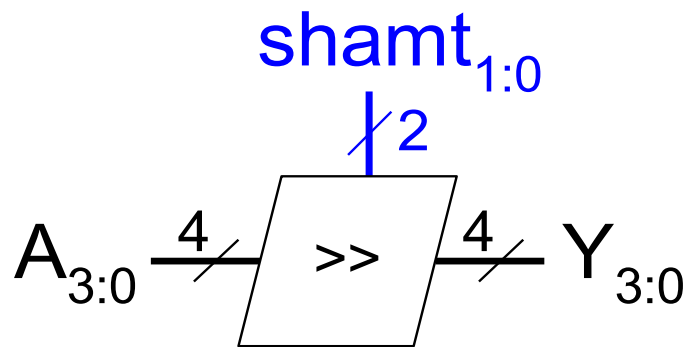
- Ex: **11001** >>> 2 =
- Ex: **11001** <<< 2 =

**Rotator:** rotates bits in a circle, such that bits shifted off one end are shifted into the other end

- Ex: **11001** ROR 2 =
- Ex: **11001** ROL 2 =



# Shifter Design



# Shifters as Multipliers, Dividers

- $A \ll N = A \times 2^N$ 
  - **Example:**  $00001 \ll 2 = 00100$  ( $1 \times 2^2 = 4$ )
  - **Example:**  $11101 \ll 2 = 10100$  ( $-3 \times 2^2 = -12$ )
- $A \gg N = A \div 2^N$ 
  - **Example:**  $01000 \gg 2 = 00010$  ( $8 \div 2^2 = 2$ )
  - **Example:**  $10000 \gg 2 = 11100$  ( $-16 \div 2^2 = -4$ )



# Multipliers

- **Partial products** formed by multiplying a single digit of the multiplier with multiplicand
- **Shifted** partial products **summed** to form result

## Decimal

$$\begin{array}{r} 230 \\ \times 42 \\ \hline 460 \\ + 920 \\ \hline 9660 \end{array}$$

$$230 \times 42 = 9660$$

## Binary

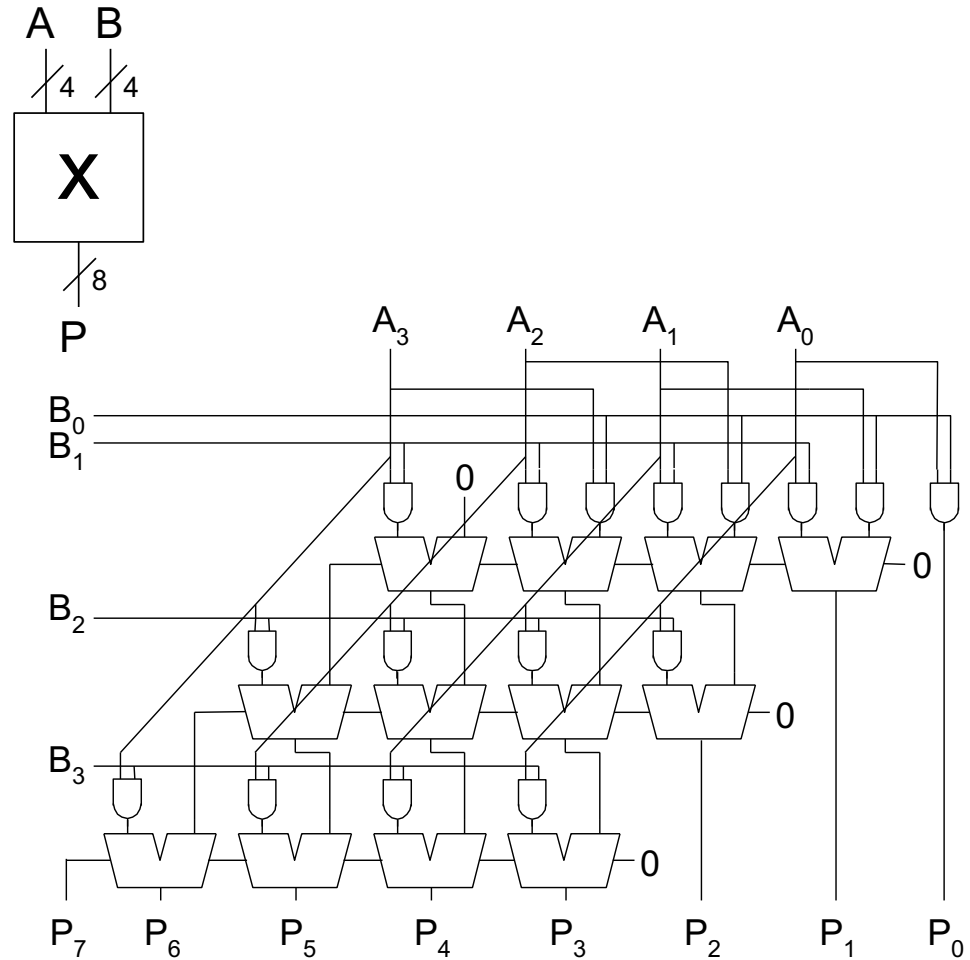
multiplicand	0101
multiplier	x 0111
partial products	<hr/> 0101 0101 0101 + 0000
result	<hr/> 0100011

$$5 \times 7 = 35$$



# 4 x 4 Multiplier

$$\begin{array}{r}
 \phantom{+} \phantom{A_3} \phantom{A_2} \phantom{A_1} \phantom{A_0} \\
 \phantom{+} \phantom{A_3} \phantom{A_2} \phantom{A_1} \phantom{A_0} \\
 \phantom{+} \phantom{A_3} \phantom{A_2} \phantom{A_1} \phantom{A_0} \\
 \phantom{+} \phantom{A_3} \phantom{A_2} \phantom{A_1} \phantom{A_0} \\
 + \phantom{A_3} \phantom{A_2} \phantom{A_1} \phantom{A_0} \\
 \hline
 P_7 \phantom{P_6} \phantom{P_5} \phantom{P_4} \phantom{P_3} \phantom{P_2} \phantom{P_1} \phantom{P_0}
 \end{array}$$



# Dividers

$$A/B = Q + R/B$$

**Decimal Example:**  $2584/15 = 172 \text{ R}4$





# Dividers

$$A/B = Q + R/B$$

**Decimal Example:**  $2584/15 = 172 \text{ R}4$

**Long-Hand:**

$$\begin{array}{r} 172 \text{ R}4 \\ 15 \overline{) 2584} \\ \underline{-15} \phantom{00} \\ 108 \phantom{0} \\ \underline{-105} \phantom{0} \\ 34 \\ \underline{-30} \\ 4 \end{array}$$



# Dividers

$$A/B = Q + R/B$$

**Decimal Example:**  $2584/15 = 172 \text{ R}4$

**Long-Hand:**

$$\begin{array}{r} 172 \text{ R}4 \\ 15 \overline{) 2584} \\ \underline{-15} \phantom{00} \\ 108 \phantom{0} \\ \underline{-105} \phantom{0} \\ 34 \\ \underline{-30} \\ 4 \end{array}$$

**Long-Hand Revisited:**

$$\begin{array}{r} 0002 \\ - 15 \\ \hline -13 \end{array} \quad \begin{array}{r} 0 \\ 3 \phantom{0} \phantom{0} \phantom{0} \\ \hline 2 \phantom{0} \phantom{0} \phantom{0} \\ 1 \phantom{0} \phantom{0} \\ 0 \end{array}$$
$$\begin{array}{r} 0025 \\ - 15 \\ \hline 10 \end{array} \quad \begin{array}{r} 0 \phantom{0} 1 \\ 3 \phantom{0} \phantom{0} \phantom{0} \\ \hline 2 \phantom{0} \phantom{0} \phantom{0} \\ 1 \phantom{0} \phantom{0} \\ 0 \end{array}$$
$$\begin{array}{r} 0108 \\ - 105 \\ \hline 3 \end{array} \quad \begin{array}{r} 0 \phantom{0} 1 \phantom{0} 7 \\ 3 \phantom{0} \phantom{0} \phantom{0} \phantom{0} \\ \hline 2 \phantom{0} \phantom{0} \phantom{0} \phantom{0} \\ 1 \phantom{0} \phantom{0} \phantom{0} \\ 0 \end{array}$$
$$\begin{array}{r} 0034 \\ - 30 \\ \hline 4 \end{array} \quad \begin{array}{r} 0 \phantom{0} 1 \phantom{0} 7 \phantom{0} 2 \\ 3 \phantom{0} \phantom{0} \phantom{0} \phantom{0} \\ \hline 2 \phantom{0} \phantom{0} \phantom{0} \phantom{0} \\ 1 \phantom{0} \phantom{0} \phantom{0} \\ 0 \end{array}$$



# Dividers

$$A/B = Q + R/B$$

**Decimal:** 2584/15 = 172 R4    **Binary:** 1101/0010 = 0110 R1

$$\begin{array}{r} 0002 \\ - 15 \\ \hline -13 \end{array}$$

$$\begin{array}{r} 0 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0001 \\ - 0010 \\ \hline 1111 \end{array}$$

$$\begin{array}{r} 0 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0025 \\ - 15 \\ \hline 10 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0011 \\ - 0010 \\ \hline 0001 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0108 \\ - 105 \\ \hline 3 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \quad 7 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0010 \\ - 0010 \\ \hline 0000 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \quad 1 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0034 \\ - 30 \\ \hline 4 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \quad 7 \quad 2 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$$

$$\begin{array}{r} 0001 \\ - 0010 \\ \hline 1111 \end{array}$$

$$\begin{array}{r} 0 \quad 1 \quad 1 \quad 0 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array} \quad R1$$



# Divider Algorithm

$$A/B = Q + R/B$$

$$R' = 0$$

for  $i = N-1$  to  $0$

$$R = \{R' \ll 1, A_i\}$$

$$D = R - B$$

if  $D < 0$ ,  $Q_i = 0$ ;  $R' = R$

else  $Q_i = 1$ ;  $R' = D$

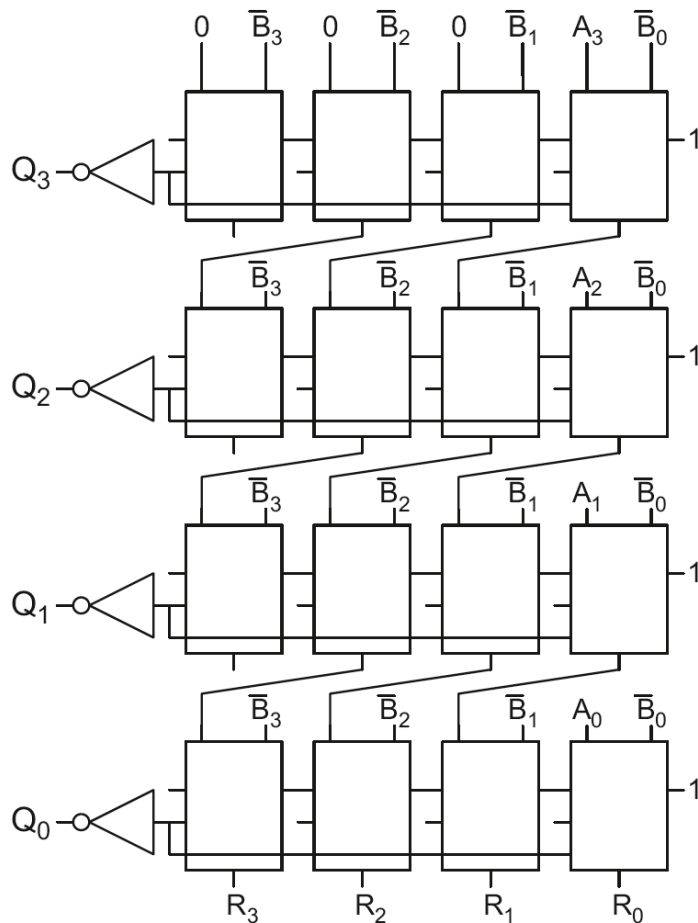
$$R = R'$$

**Binary: 1101/10 = 0110 R1**

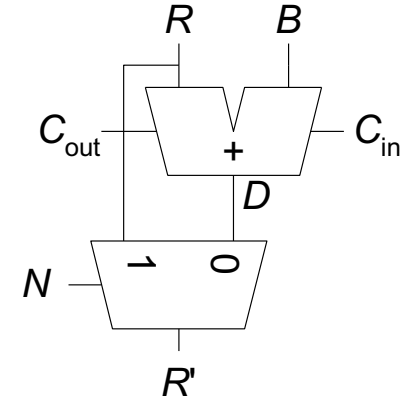
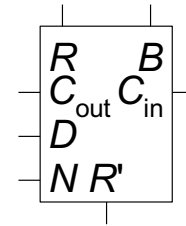
$\begin{array}{r} 0001 \\ -0010 \\ \hline 1111 \end{array}$	$\begin{array}{r} 0 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$
$\begin{array}{r} 0011 \\ -0010 \\ \hline 0001 \end{array}$	$\begin{array}{r} 0 \quad 1 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$
$\begin{array}{r} 0010 \\ -0010 \\ \hline 0000 \end{array}$	$\begin{array}{r} 0 \quad 1 \quad 1 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array}$
$\begin{array}{r} 0001 \\ -0010 \\ \hline 1111 \end{array}$	$\begin{array}{r} 0 \quad 1 \quad 1 \quad 0 \\ \hline 3 \quad 2 \quad 1 \quad 0 \end{array} \quad R1$



# 4 x 4 Divider



Legend



**Division:**  $A/B = Q + R/B$

$R' = 0$

for  $i = N-1$  to  $0$

$R = \{R' \ll 1, A_i\}$

$D = R - B$

if  $D < 0$ ,  $Q_i = 0$ ,  $R' = R$

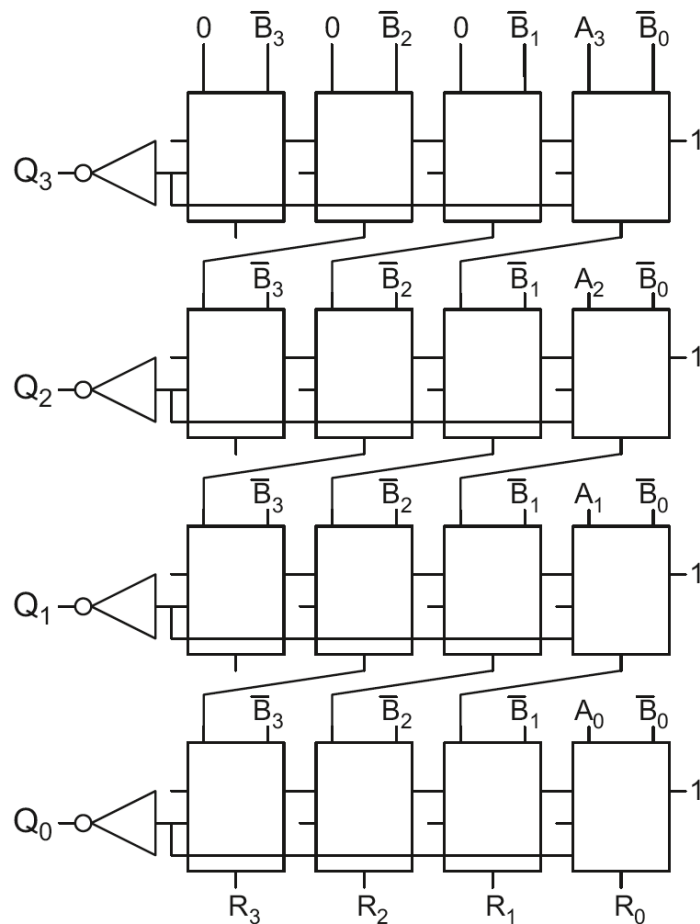
else  $Q_i = 1$ ,  $R' = D$

$R = R'$

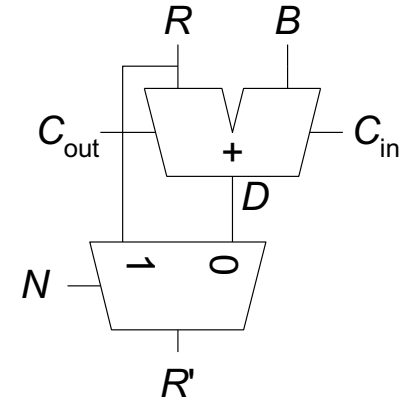
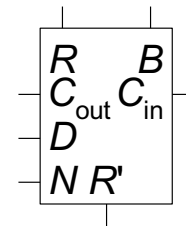
Each row computes one iteration of the division algorithm.



# 4 x 4 Divider



Legend



$$\begin{array}{r} 0001 \\ -0010 \\ \hline 1111 \end{array} \quad \begin{array}{r} 0 \\ \hline 3 \ 2 \ 1 \ 0 \end{array}$$

$$\begin{array}{r} 0011 \\ -0010 \\ \hline 0001 \end{array} \quad \begin{array}{r} 0 \ 1 \\ \hline 3 \ 2 \ 1 \ 0 \end{array}$$

$$\begin{array}{r} 0010 \\ -0010 \\ \hline 0000 \end{array} \quad \begin{array}{r} 0 \ 1 \ 1 \\ \hline 3 \ 2 \ 1 \ 0 \end{array}$$

$$\begin{array}{r} 0001 \\ -0010 \\ \hline 1111 \end{array} \quad \begin{array}{r} 0 \ 1 \ 1 \ 0 \\ \hline 3 \ 2 \ 1 \ 0 \end{array} \quad R1$$

Each row computes one iteration of the division algorithm.



# Number Systems

Numbers we can represent using binary representations

- **Positive numbers**
  - Unsigned binary
- **Negative numbers**
  - Two's complement
  - Sign/magnitude numbers

What about **fractions**?



# Numbers with Fractions

Two common notations:

- **Fixed-point:** binary point fixed
- **Floating-point:** binary point floats to the right of the most significant 1





# Fixed-Point Numbers

- 6.75 using 4 integer bits and 4 fraction bits:

01101100

0110.1100

$$2^2 + 2^1 + 2^{-1} + 2^{-2} = 6.75$$

- Binary point is implied
- The number of integer and fraction bits must be agreed upon beforehand



# Fixed-Point Number Example

- Represent  $7.5_{10}$  using 4 integer bits and 4 fraction bits.



# Signed Fixed-Point Numbers

- **Representations:**

- Sign/magnitude
- Two's complement

- **Example:** Represent  $-7.5_{10}$  using 4 integer and 4 fraction bits

- **Sign/magnitude:**

- **Two's complement:**

1. +7.5:

2. Invert bits:

3. Add 1 to lsb:  $\quad + \quad \quad \quad \underline{\quad 1 \quad}$



# Floating-Point Numbers

- Binary point floats to the right of the most significant 1
- Similar to decimal scientific notation

- For example, write  $273_{10}$  in scientific notation:

$$273 = 2.73 \times 10^2$$

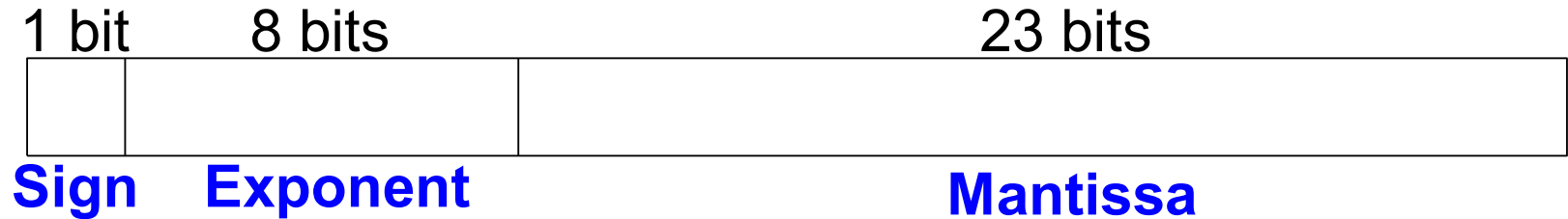
- In general, a number is written in scientific notation as:

$$\pm M \times B^E$$

- **M** = mantissa
- **B** = base
- **E** = exponent
- In the example,  $M = 2.73$ ,  $B = 10$ , and  $E = 2$



# Floating-Point Numbers



- **Example:** represent the value  $228_{10}$  using a 32-bit floating point representation

We show three versions – final version is called the **IEEE 754 floating-point standard**



# Floating-Point Representation 1

1. Convert decimal to binary

$$228_{10} = 11100100_2$$

2. Write the number in “binary scientific notation”:

$$11100100_2 = 1.11001_2 \times 2^7$$

3. Fill in each field of the 32-bit floating point number:

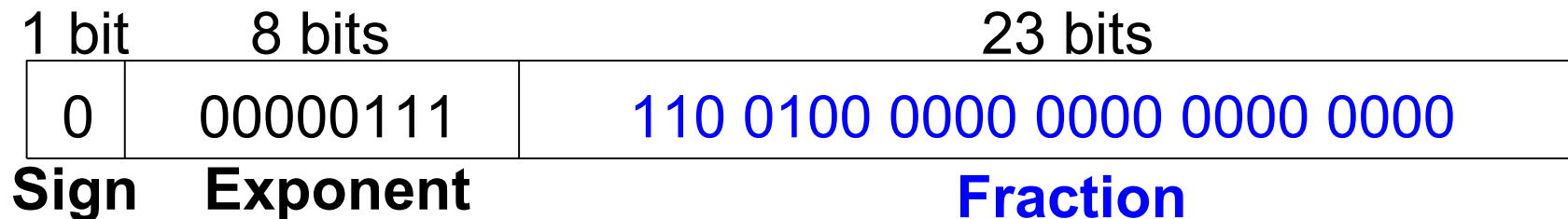
- The sign bit is positive (0)
- The 8 exponent bits represent the value 7
- The remaining 23 bits are the mantissa

1 bit	8 bits	23 bits
0	00000111	11 1001 0000 0000 0000 0000
<b>Sign</b>	<b>Exponent</b>	<b>Mantissa</b>



# Floating-Point Representation 2

- First bit of the mantissa is always 1:
  - $228_{10} = 11100100_2 = \mathbf{1.11001} \times 2^7$
- So, no need to store it: *implicit leading 1*
- Store just fraction bits in 23-bit field



# Floating-Point Representation 3

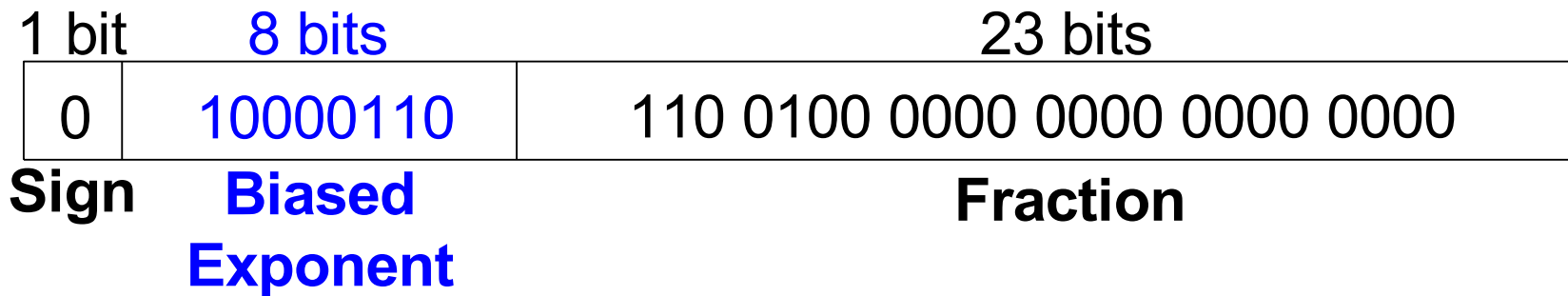
- *Biased exponent*: bias = 127 ( $01111111_2$ )

- Biased exponent = bias + exponent

- Exponent of 7 is stored as:

$$127 + 7 = 134 = 0x10000110_2$$

- The **IEEE 754 32-bit floating-point representation** of  $228_{10}$



in hexadecimal: **0x43640000**





# Floating-Point Example

Write  $-58.25_{10}$  in floating point (IEEE 754)



# Floating-Point Example

Write  $-58.25_{10}$  in floating point (IEEE 754)

1. Convert decimal to binary:

$$58.25_{10} = \mathbf{111010.01}_2$$

2. Write in binary scientific notation:

$$\mathbf{1.1101001} \times 2^5$$

3. Fill in fields:

**Sign bit:** **1** (negative)

**8 exponent bits:**  $(127 + 5) = 132 = \mathbf{10000100}_2$

**23 fraction bits:**  $\mathbf{110\ 1001\ 0000\ 0000\ 0000\ 0000}$

1 bit	8 bits	23 bits
1	100 0010 0	110 1001 0000 0000 0000 0000
<b>Sign</b>	<b>Exponent</b>	<b>Fraction</b>

in hexadecimal:  $\mathbf{0xC2690000}$



# Floating-Point: Special Cases

Number	Sign	Exponent	Fraction
0	X	00000000	00000000000000000000000000000000
$\infty$	0	11111111	00000000000000000000000000000000
$-\infty$	1	11111111	00000000000000000000000000000000
NaN	X	11111111	non-zero



# Floating-Point Precision

- **Single-Precision:**
  - 32-bit
  - 1 sign bit, 8 exponent bits, 23 fraction bits
  - bias = 127
  
- **Double-Precision:**
  - 64-bit
  - 1 sign bit, 11 exponent bits, 52 fraction bits
  - bias = 1023



# Floating-Point: Rounding

- **Overflow:** number too large to be represented
- **Underflow:** number too small to be represented
- **Rounding modes:**
  - Down
  - Up
  - Toward zero
  - To nearest
- **Example:** round 1.100101 (1.578125) to only 3 fraction bits
  - Down: 1.100
  - Up: 1.101
  - Toward zero: 1.100
  - To nearest: 1.101 (1.625 is closer to 1.578125 than 1.5 is)



# Floating-Point Addition

1. Extract exponent and fraction bits
2. Prepend leading 1 to form mantissa
3. Compare exponents
4. Shift smaller mantissa if necessary
5. Add mantissas
6. Normalize mantissa and adjust exponent if necessary
7. Round result
8. Assemble exponent and fraction back into floating-point format



# Floating-Point Addition Example

Add the following floating-point numbers:

0x3FC00000

0x40500000



# Floating-Point Addition Example

## 1. Extract exponent and fraction bits

1 bit	8 bits	23 bits
0	01111111	100 0000 0000 0000 0000 0000
<b>Sign</b>	<b>Exponent</b>	<b>Fraction</b>

1 bit	8 bits	23 bits
0	10000000	101 0000 0000 0000 0000 0000
<b>Sign</b>	<b>Exponent</b>	<b>Fraction</b>

For first number (N1):  $S = 0, E = 127, F = .1$

For second number (N2):  $S = 0, E = 128, F = .101$

## 2. Prepend leading 1 to form mantissa

N1: 1.1

N2: 1.101





# Floating-Point Addition Example

## 3. Compare exponents

$127 - 128 = -1$ , so shift N1 right by 1 bit

## 4. Shift smaller mantissa if necessary

shift N1's mantissa:  $1.1 \gg 1 = 0.11$  ( $\times 2^1$ )

## 5. Add mantissas

$$\begin{array}{r} 0.11 \times 2^1 \\ + 1.101 \times 2^1 \\ \hline 10.011 \times 2^1 \end{array}$$



# Floating Point Addition Example

## 6. Normalize mantissa and adjust exponent if necessary

$$10.011 \times 2^1 = 1.0011 \times 2^2$$

## 7. Round result

No need (fits in 23 bits)

## 8. Assemble exponent and fraction back into floating-point format

$$S = 0, E = 2 + 127 = 129 = 10000001_2, F = 001100..$$

1 bit	8 bits	23 bits
0	10000001	001 1000 0000 0000 0000 0000
<b>Sign</b>	<b>Exponent</b>	<b>Fraction</b>

in hexadecimal: **0x40980000**

