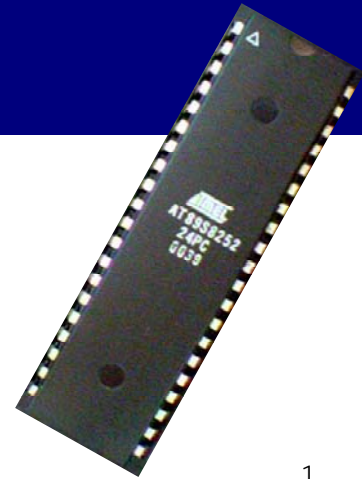# Digital Integrated Circuits & Microcontrollers

**Chapter 2. Binary representation**

---

## Numeration Systems

- ## Decimal
  - □ the set of symbols is:
    {0, 1, 2, 3, 4, 5, 6, 7, 8, 9};
  - □ The integer number 493 in base 10:
    $493_{10} = 4 \cdot 10^2 + 9 \cdot 10^1 + 3 \cdot 10^0 =$
    $= 400 \quad + \quad 90 \quad + 3$
  - □ The *rational* number 35,54 in base 10:
    $35.64_{10} = 3 \cdot 10^1 + 5 \cdot 10^0 + 6 \cdot 10^{-1} + 4 \cdot 10^{-2} =$
    $= 30 \quad + \quad 5 \quad + \quad 0.6 \quad + \quad 0.04$

## Numeration Systems

- ## Binary
  - ☐ the set of symbols is:
    {0, 1};
  - ☐ $11001_2 = 1 \cdot 2^4 + 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 =$
    $= 16 + 8 + 1 = 25$
  - ☐ $110.01_2 = 1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} =$
    $= 4 + 2 + 0.25 = 6.25$

## Conversion between numeration systems

- ## Conversion from binary into hexadecimal
  - ☐ Group into 4 bits (nibbles)
  - ☐ Each nibble corresponds to a hexadecimal symbol:

$$1101101,1001101_2$$

$$0110 \quad 1101,1001 \quad 1010_2$$

$$6 \quad\quad D \quad , \quad 9 \quad\quad A_{16}$$

# Conversion between numeration systems

- ## Conversion from decimal into base *b*
  - ☐ For the **integer part** we can write:

  $$i = \left(\left(\left(\left(d_k b + d_{k-1}\right)b + ... + d_3\right)b + d_2\right)b + d_1\right)b + d_0$$

    - $d_0$ is the remainder of division of *i* to **b**
    - The quotient (another integer) is divided to **b**
    - Repeat until reach 0.
    - The remainder obtained after each division is the symbol $d_k$ of representing into the base **b**.

  $$i = d_k b^k + d_{k-1} b^{k-1} + ... + d_3 b^3 + d_2 b^2 + d_1 b^1 + d_0 b^0$$

# Conversion between numeration systems

- ## Conversion from decimal into base *b*
  - ☐ For the **fractional part**:

  $$f = d_{-1} b^{-1} + d_{-2} b^{-2} + d_{-3} b^{-3} + ... + d_{-k} b^{-k} + ...$$

    - Multiply *f* with b

  $$b \cdot f = d_{-1} + d_{-2} b^{-1} + d_{-3} b^{-2} + ... + d_{-k} b^{-k+1} + ...$$

    - Keep the integer part from the right part, $d_{-1}$, which is subtracted from the left part.
    - Continue by multiplying the remaining fractionary part to b until reach 0.

  $$b\left(b \cdot f - d_{-1}\right) = d_{-2} + d_{-3} b^{-1} + ... + d_{-k} b^{-k+2} + ...$$

## Conversion between numeration systems

- **EXAMPLE:**
  - Represent the number 23.65 into base 2:
    - The integer part:

$$23 : 2 = 11 \quad | \quad 1 \quad \text{LSB}$$
$$11 : 2 = 5 \quad | \quad 1$$
$$5 : 2 = 2 \quad | \quad 1$$
$$2 : 2 = 1 \quad | \quad 0$$
$$1 : 2 = 0 \quad | \quad 1 \quad \text{MSB}$$

$$23_{10} = \mathbf{10111}_2$$

## Conversion between numeration systems

- The fractional part:

$$0.65 \times 2 = 1.3 \quad | \quad 1 \quad \text{MSB}$$
$$0.3 \ \times 2 = 0.6 \quad | \quad 0$$
$$0.6 \ \times 2 = 1.2 \quad | \quad 1$$
$$0.2 \ \times 2 = 0.4 \quad | \quad 0$$
$$0.4 \ \times 2 = 0.8 \quad | \quad 0 \qquad 0.65_{10} = \mathbf{0.10(1001)}_2$$
$$0.8 \ \times 2 = 1.6 \quad | \quad 1$$
$$0.6 \ \times 2 = 1.2 \quad | \quad 1$$
$$0.2 \ \times 2 = 0.4 \quad | \quad 0 \quad \dots$$

- Results that number 0.65 cannot be exactly represented on a finite number of bits.
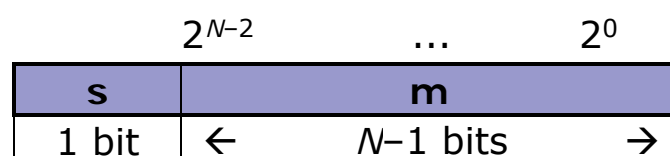
## Negative numbers representation

- **MSB – sign bit.**
  - □ 0 for positive numbers (+);
  - □ 1 for negative numbers (–).
- **The rest of $N-1$ bits are for value representation.**

| | $2^{N-2}$ ... $2^0$ |
|---|---|
| **s** | **m** |
| 1 bit | ← $N-1$ bits → |

---

## Negative numbers representation

- **The representation: sign bit, magnitude**
  - □ **EXAMPLE:**
    - 9 = 0 01001
    - –9 = 1 01001
  - □ The range of representation:
    - $2^{N-1}$ positive values between 0 and $2^{N-1}-1$.
    - $2^{N-1}$ negative values between $-(2^{N-1}-1)$ and 0.

| | $2^{N-2}$ ... $2^0$ |
|---|---|
| **s** | **m** |
| 1 bit | ← $N-1$ bits → |

# Negative numbers representation

- **Two's complement** representation
  - ☐ The negative numbers representation is obtained by addition of $2^N$.
  - ☐ **EXAMPLE:**
    - For $N$=6 bits ($2^N = 64$).

      $13 = 001101_2$

      $-13$ corresponds to $64 + (-13) = 51 = 110011_2$

---

# Two's complement representation

- For obtaining negative numbers:
  - Each bit is complemented;
  - Add 1.
  - ☐ **EXAMPLE:**
    - for                       $13 = 001101_2$

      complement each bit:   $110010_2$

      add 1:                     $110011_2 = -13$
    - for                       $-13 = 110011_2$

      complement each bit:   $001100_2$
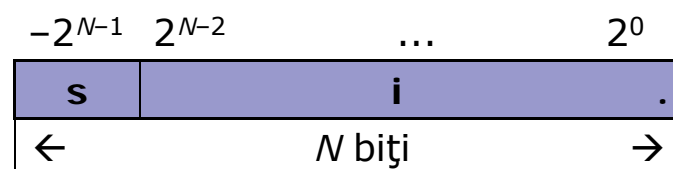
      add 1:                     $001101_2 = +13$

## Two's complement representation

- The range of two's complement representation
  - $2^{N-1}$ positive values between 0 and $2^{N-1}-1$.
  - $2^{N-1}$ negative values between $-2^{N-1}$ and -1.
- The result of adding a number with its two's complement is 0:

| | |
|---:|:---:|
| 13+ | $001101_2$ |
| $-13$ | $110011_2$ |
| = 0 | $1000000_2$ |

## Integer numbers representation

- Binary representation is considered **right aligned** (decimal point is right to LSB).
- MSB represents the sign bit.



- Integer two's complement range:

$$-2^{N-1},...,-1,0,...,2^{N-1}-1$$

## Integer numbers representation

- **EXAMPLE:** (for *N*=4 bits)

| Decimal | Binary |
|---------|--------|
| 0 | 0000 |
| 1 | 0001 |
| 2 | 0010 |
| 3 | 0011 |
| 4 | 0100 |
| 5 | 0101 |
| 6 | 0110 |
| 7 | 0111 |

| Decimal | Binary |
|---------|--------|
| -8 | 1000 |
| -7 | 1001 |
| -6 | 1010 |
| -5 | 1011 |
| -4 | 1100 |
| -3 | 1101 |
| -2 | 1110 |
| -1 | 1111 |

## Integer numbers representation

- *Sign integers addition*

```
    1111              011              11
+3 | 0011       -5 | 1011       -3 | 1101
-2 | 1110       +3 | 0011       -4 | 1100
 1 | 0001       -2 | 1110       -7 | 1001
```

  □ An overflow occurs if the result is outside of the *N* bits representation range:

```
+3 | 0011              -3 | 1101
+6 | 0110              -6 | 1010
 9 | 1001=-7           -9 | 0111 =7
```

## Integer numbers representation

- ***Sign bit extension***
  - ☐ Needed when **increasing** the number of bits for the integer part.
  - ☐ Sign bit is copied to the left toward MSB.

|  | $N=4$ biţi |
|---|---|
| +3 | **0**011 |
| −3 | **1**101 |

|  | $N'=8$ biţi |
|---|---|
| +3 | **0000 0**011 |
| −3 | **1111 1**101 |

## Integer numbers representation

- ***Multiplying by a power of 2***
  - ☐ **Multiplying** by $2^k$ is equivalent with **shifting left $k$ bits** and **filling with 0** toward LSB.

|  | $N=8$ biţi |
|---|---|
| 3 | 0000 0011 |
| $3{\cdot}2^2$ | 0000 11**00** |

|  | $N=8$ biţi |
|---|---|
| −3 | 1111 1101 |
| $-3{\cdot}2^2$ | 1111 01**00** |

## Integer numbers representation

- **_Dividing by a power of 2_**
  - □ **Dividing** by $2^k$ is equivalent with **shifting right _k_ bits** and **sign bit extension**.

|            | $N$=8 biţi |
| ---------- | --------- |
| 24         | 0001 1000 |
| $24/2^3$   | **000**0 0011 |

|              | $N$=8 biţi |
| ------------ | --------- |
| $-24$        | 1110 1000 |
| $-24/2^3$    | **111**1 1101 |

## Integer numbers representation

- **_Integers multiplication_**
  - □ **unsigned integers**
  - □ The result in double precision representation
  - □ Binary multiplication with 0 and 1

```
  6 |    0110
 ×5 |   ×0101
    |    0110
    |   0000
    |   0110
 ___|  0000      +
 30 | 0011110
```

## Integer numbers representation

- **Multiplication is equivalent with consecutive shift and add operations.**
  - For example 5 can be expressed:
    $$5 = 2^0 + 2^2$$
  - Multiplication can be computed:
    $$6 \times 5 = 6 \times \left(2^0 + 2^2\right) = 6 \times 1 + 6 \times 2^2$$

$$
\begin{array}{r|l}
6 \times 1 & \phantom{00}0110 \\
6 \times 2^2 & \underline{0110 \phantom{00} +} \\
& 0011110
\end{array}
$$

## Integer numbers representation

- ***Integers division***
  - Successive subtractions of the divisor from the dividend.

$$
\begin{array}{ll}
15|\underline{\phantom{0}3} & \quad \underline{1111} \;\;\big| 11 \\
\phantom{15|}5 & \quad -\underline{11} \;\;\;\big| 101 \\
& \quad \phantom{-}0011 \\
& \quad \phantom{00}-\underline{11} \\
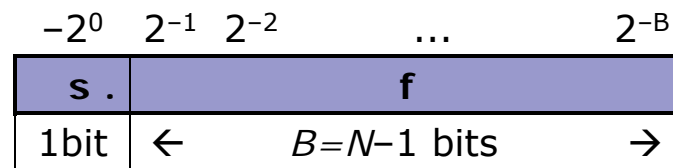& \quad \phantom{000}00 \;\;\big|
\end{array}
$$

## Fractional numbers representation

- A fractional part *f* is any number who's modulus satisfies the inequality:

$$0.0 \leq | f | < 1.0$$

  - Left aligned: binary point is at the right of MSB

| | $-2^0$ | $2^{-1}$ | $2^{-2}$ | ... | $2^{-B}$ |
|---|---|---|---|---|---|
| **s .** | | | **f** | | |
| 1bit | $\leftarrow$ | | $B=N-1$ bits | | $\rightarrow$ |

  - Fixed point fractional representation range:

$$-1,...,-2^{-B},0,2^{-B}...,1-2^{-B}$$

## Fractional numbers representation

- **EXAMPLE:**
  - for *N*=4 bits

| Decimal | Binary |
|---|---|
| 0 | 0000 |
| 0.125 | 0001 |
| 0.250 | 0010 |
| 0.375 | 0011 |
| 0.500 | 0100 |
| 0.625 | 0101 |
| 0.750 | 0110 |
| 0.875 | 0111 |

| Decimal | Binary |
|---|---|
| -1 | 1000 |
| -0.875 | 1001 |
| -0.750 | 1010 |
| -0.625 | 1011 |
| -0.500 | 1100 |
| -0.375 | 1101 |
| -0.250 | 1110 |
| -0.125 | 1111 |

## Fractional numbers representation

- **Q*m.n* format**
  - □ *n* bits for fractional part;
  - □ (optional) specify the number of bits *m* for the integer part, excluding the sign bit (MSB);
  - □ The complete binary representation has *1+m+n* bits.
- **EXAMPLE** (for *N*=16 bits):
  - □ Q15 means 15 bits for the fractional part (16 bits with the sign bit)
  - □ Q1.14 has 1 bit for the integer part, 14 bits for the fractional part and the sign bit.

---

## Fractional numbers representation

- *Quick conversion of fractional numbers into binary*
  - □ *f* represented on *N*=*B*+1 bits is an **integer multiple of $2^{-B}$**

$$-2^0 \quad 2^{-1} \; 2^{-2} \qquad \ldots \qquad 2^{-B}$$
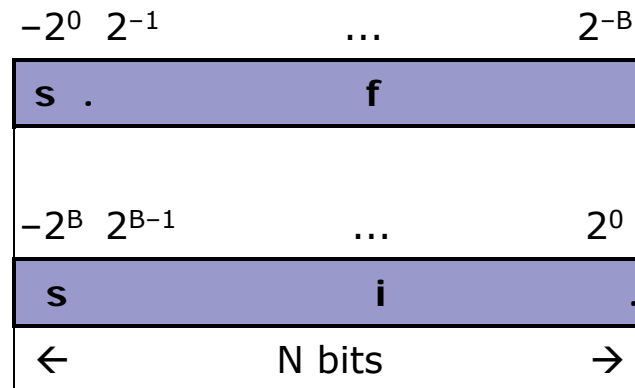
| s . | f |
|-----|---|

  - □ Let *i* the corresponding integer multiple:

$$i = f \cdot 2^B$$

    - ■ Equivalent with a left shifting of *f* with B bits.

## Fractional numbers representation



| | $-2^0$ | $2^{-1}$ | ... | $2^{-B}$ |
|---|---|---|---|---|
| s | . | | f | |

| | $-2^B$ | $2^{B-1}$ | ... | $2^0$ |
|---|---|---|---|---|
| s | | i | | . |

← N bits →

□ Converting *f* from binary to decimal
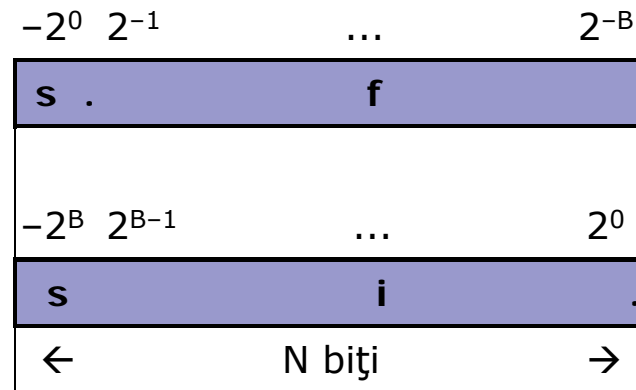- consider the binary representation for the corresponding signed integer *i* and then, divide by $2^B$

---

## Fractional numbers representation

- **EXAMPLE:**
  □ Binary to decimal (*N*=8 bits, *B*=7),
  - For:  0.010 0110
    decimal integer:  0010 0110 = $38_{10}$
    divide by $2^7$=128:  38/128 = 0.296875

  - For:  1.110 1100
    two's complement:  –0001 0100 = $-20_{10}$
    divide by $2^7$=128:  –20/128 = –0.15625

# Fractional numbers representation

$$-2^0 \quad 2^{-1} \qquad\qquad \ldots \qquad\qquad 2^{-B}$$

| s  . | f |
|------|---|
|      |   |

$$-2^B \quad 2^{B-1} \qquad\qquad \ldots \qquad\qquad 2^0$$

| s | i | . |
|---|---|---|

$$\leftarrow \qquad\qquad N \text{ biţi} \qquad\qquad \rightarrow$$

☐ Converting $f$ from decimal to binary

- multiply $f$ by $2^B$ (shift left B bits)
- Represent the integer part $i$ as a signed integer on $N$ bits.

---

# Fractional numbers representation

- ## EXAMPLE:

  ☐ Decimal to binary ($N$=8 biţi, $B$=7)

  - For:            0.875
    multiply by $2^7$:        $0.875 \cdot 128 = 112_{10}$
    represented binary:     $112_{10} = 01110000_2$

  - For:            −0.625
    multiply by $2^7$:        $-0.625 \cdot 128 = -80_{10}$
    represented in        $80_{10} = 01010000_2$
    two's complement:     $-80_{10} = 10110000_2$

# Fractional numbers representation

- For 0.65 ($N$=8 bits, $B$=7)
  multiply by $2^7$: $0.65 \cdot 128 = 83.2_{10}$
  take the integer part: $83_{10} = 01010011_2$

- In the last example the **quantization error** (by truncation) appears since the fractional 0.65 cannot be exactly represented on 8 bits.
  - The result $0.1010011_2$ is equal to 0.6484375.
  - The quantization error is:

    $$0.65$$
    $$- 0.6484375$$
    $$= 0.0015625$$

# Floating point representation
- The following are equivalent representations of 1,234

$$123{,}400.0 \quad \times 10^{-2}$$
$$12{,}340.0 \quad \times 10^{-1}$$
$$1{,}234.0 \quad \times 10^{0}$$
$$123.4 \quad \times 10^{1}$$
$$12.34 \quad \times 10^{2}$$
$$1.234 \quad \times 10^{3}$$
$$0.1234 \times 10^{4}$$

The representations differ in that the decimal place – the "point" – "floats" to the left or right (with the appropriate adjustment in the exponent).

## Parts of a Floating Point Number

$$-0.9876 \times 10^{-3}$$

Sign of mantissa

Location of decimal point
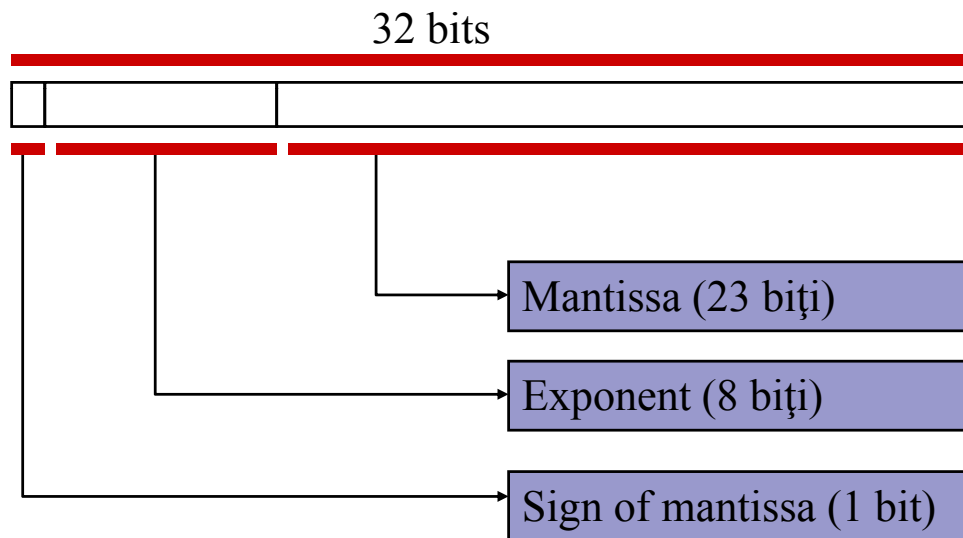
Mantissa

Sign of exponent

Exponent

Base

## IEEE 754 Standard

- Single precision: 32 bits, consisting of...
  - Sign bit (1 bit)
  - Exponent (8 bits)
  - Mantissa (23 bits)
- Double precision: 64 bits, consisting of…
  - Sign bit (1 bit)
  - Exponent (11 bits)
  - Mantissa (52 bits)

## Single Precision Format

Mantissa (23 biți)

Exponent (8 biți)

Sign of mantissa (1 bit)

## Normalization

- The mantissa is *normalized*
  - Has an implied decimal place on left
  - Has an implied "1" on left of the decimal place
- E.g.,
  - Mantissa $\rightarrow$ 10100000000000000000000
  - Represents... $1.101_2 = 1.625_{10}$
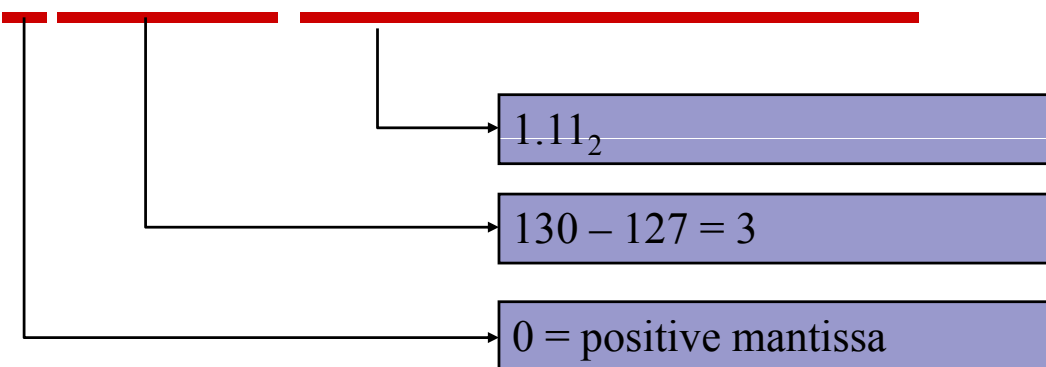
## Excess Notation

- To include +ve and –ve exponents, "excess" notation is used
  - □ Single precision:  excess 127
  - □ Double precision: excess 1023
- The value of the exponent stored is larger than the actual exponent
- E.g., excess 127,
  - □ Exponent $\rightarrow$     `10000111`
  - □ Represents…   `135 - 127 = 8`

---

## Example

- Single precision

`0 10000010 11000000000000000000000`

$1.11_2$

$130 - 127 = 3$

$0 =$ positive mantissa

$+1.11_2 \times 2^3 = 1110.0_2 = 14.0_{10}$

## Hexadecimal

■ It is convenient and common to represent the original floating point number in hexadecimal

■ The preceding example…

```
0 10000010 11000000000000000000000
```

$$4 \quad 1 \quad 6 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

## Converting from Floating Point

■ What decimal value is represented by the following 32-bit floating point number?

$$C17B0000_{16}$$

# Converting _from_ Floating Point

- Step 1
  - □ Express in binary and find S, E, and M

    $C17B0000_{16}$ =

    <u>1</u>  <u>10000010</u>  <u>11110110000000000000000</u>$_2$
    S      E                    M

    1 = negative
    0 = positive

---

# Converting _from_ Floating Point

- Step 2
  - □ Find "real" exponent, _n_
  - □ _n_ = E – 127
    - = $10000010_2$ – 127
    - = 130 – 127
    - = 3

## Converting from Floating Point

- Step 3
  - □ Put S, M, and *n* together to form binary result
  - □ (Don't forget the implied "1." on the left of the mantissa.)

$$-1.1111011_2 \times 2^n =$$

$$-1.1111011_2 \times 2^3 =$$

$$-1111.1011_2$$

## Converting from Floating Point

- Step 4
  - □ Express result in decimal

$$\underline{-1111}.\underline{1011}_2$$

-15

$2^{-1} = 0.5$
$2^{-3} = 0.125$
$2^{-4} = \underline{0.0625}$
$\phantom{2^{-4} = }0.6875$

## Answer: -15.6875

## Converting to Floating Point

- Express $36.5625_{10}$ as a 32-bit floating point number (in hexadecimal)

## Converting to Floating Point

- Step 1
  - □ Express original value in binary

    $36.5625_{10}$ =

    $100100.1001_{2}$

## Converting to Floating Point
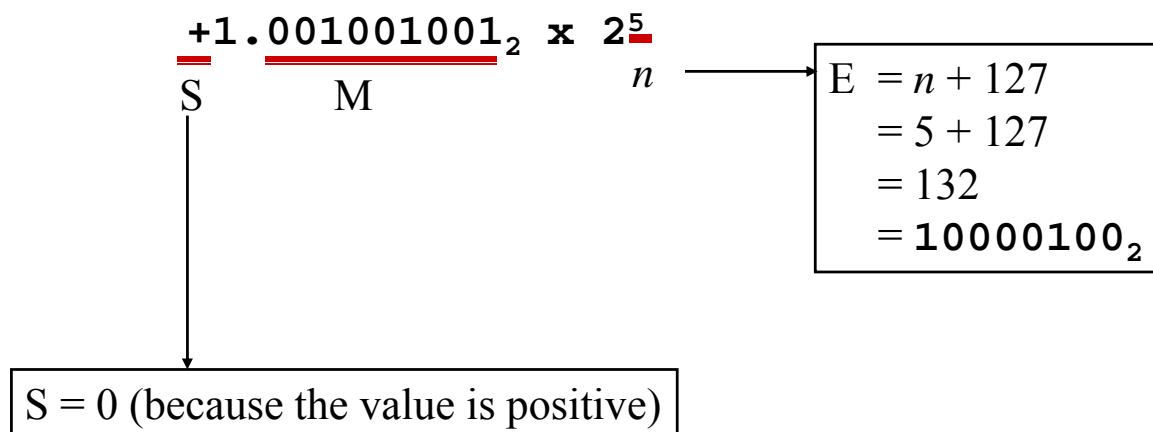
- Step 2
  - Normalize

$$100100.1001_2 =$$

$$1.001001001_2 \times 2^5$$

## Converting to Floating Point

- Step 3
  - Determine S, E, and M

$$\underset{\text{S}}{+}\underset{\text{M}}{1.001001001}_2 \times 2^{\underset{n}{5}}$$

$$
\begin{aligned}
E &= n + 127 \\
&= 5 + 127 \\
&= 132 \\
&= 10000100_2
\end{aligned}
$$

S = 0 (because the value is positive)

## Converting to Floating Point

- Step 4
  - Put S, E, and M together to form 32-bit binary result

$$0 \quad 10000100 \quad 00100100100000000000000_2$$

$$\text{S} \qquad \text{E} \qquad\qquad \text{M}$$

## Converting to Floating Point

- Step 5
  - Express in hexadecimal

```
0 10000100 00100100100000000000000₂ =

0100 0010 0001 0010 0100 0000 0000 0000₂ =

  4    2    1    2    4    0    0    0₁₆
```

$$\text{Answer: } 42124000_{16}$$