# + Chapter 10 Computer Arithmetic

STUDENTS-HUB.com

# Arithmetic - The heart of instruction execution



# Arithmetic & Logic Unit (ALU)

- Part of the computer that actually performs arithmetic and logical operations on data
- All of the other elements of the computer system are there mainly to bring data into the ALU for it to process and then to take the results back out
- Based on the use of simple digital logic devices that can store binary digits and perform simple Boolean logic operations
- Does the calculations
- Handles integers
- May handle floating point (real) numbers
- May be separate FPU (maths coprocessor).
- May be on chip separate FPU (486DX +)

3



#### Figure 10.1 ALU Inputs and Outputs

# **Positional Number Systems**

Different Representations of Natural Numbers

- XXVII Roman numerals (not positional)
- 27 Radix-10 or decimal number (positional)
- 11011<sub>2</sub> Radix-2 or binary number (also positional)

Fixed-radix positional representation with k digits

Number N in radix 
$$r = (d_{k-1}d_{k-2} \dots d_1d_0)_r$$
  
Value =  $d_{k-1} \times r^{k-1} + d_{k-2} \times r^{k-2} + \dots + d_1 \times r + d_0$   
Examples:  $(11011)_2 = 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2 + 1 = 27$   
 $(2103)_4 = 2 \times 4^3 + 1 \times 4^2 + 0 \times 4 + 3 = 147$ 

5

### **Binary Numbers**

- Each binary digit (called bit) is either 1 or 0
- Bits have no inherent meaning, can represent
  - Unsigned and signed integers
  - Characters
  - Floating-point numbers
  - Images, sound, etc.
- Bit Numbering



- Least significant bit (LSB) is rightmost (bit 0)
- Most significant bit (MSB) is leftmost (bit 7 in a 8-bit number)

6

# **Hexadecimal Integers**

- 16 Hexadecimal Digits: 0 9, A F
- More convenient to use than binary numbers

Binary, Decimal, and Hexadecimal Equivalents

| Binary | Decimal | Hexadecimal | Binary | Decimal | Hexadecimal |
|--------|---------|-------------|--------|---------|-------------|
| 0000   | 0       | 0           | 1000   | 8       | 8           |
| 0001   | 1       | 1           | 1001   | 9       | 9           |
| 0010   | 2       | 2           | 1010   | 10      | А           |
| 0011   | 3       | 3           | 1011   | 11      | В           |
| 0100   | 4       | 4           | 1100   | 12      | С           |
| 0101   | 5       | 5           | 1101   | 13      | D           |
| 0110   | 6       | 6           | 1110   | 14      | Е           |
| 0111   | 7       | 7           | 1111   | 15      | F           |

STUDENTS-HUB.com

# **Converting Binary to Hexadecimal**

Each hexadecimal digit corresponds to 4 binary bits

Example:

Convert the 32-bit binary number to hexadecimal

1110 1011 0001 0110 1010 0111 1001 0100

Solution:

| E    | в    | 1    | 6    | A    | 7    | 9    | 4    |
|------|------|------|------|------|------|------|------|
| 1110 | 1011 | 0001 | 0110 | 1010 | 0111 | 1001 | 0100 |

8

## **Integer Storage Sizes**

| Byte [<br>Half Word [ | 8  | Storage | e Sizes |
|-----------------------|----|---------|---------|
| Word                  | 32 |         |         |
| Double Word           |    | 64      |         |

| Storage Type | Unsigned Range                  | Powers of 2                |
|--------------|---------------------------------|----------------------------|
| Byte         | 0 to 255                        | 0 to (2 <sup>8</sup> – 1)  |
| Half Word    | 0 to 65,535                     | 0 to (2 <sup>16</sup> – 1) |
| Word         | 0 to 4,294,967,295              | 0 to (2 <sup>32</sup> – 1) |
| Double Word  | 0 to 18,446,744,073,709,551,615 | 0 to (2 <sup>64</sup> – 1) |

What is the largest 20-bit unsigned integer?

Answer:  $2^{20} - 1 = 1,048,575$ 

### **Signed Integers**

Several ways to represent a signed number

- Sign-Magnitude
- Biased
- l's complement
- 2's complement

Divide the range of values into 2 equal parts

- First part corresponds to the positive numbers  $(\geq 0)$
- Second part correspond to the negative numbers (< 0)</p>

Focus will be on the 2's complement representation

- Has many advantages over other representations
- Used widely in processors to represent signed integers

| Decimal | Signed<br>Magnitude | Ones<br>Complement | Twos<br>Complement | Biased<br>B=+8 | Biased<br>B=+7 |
|---------|---------------------|--------------------|--------------------|----------------|----------------|
| +8      | -                   | -                  | -                  | -              | 1111           |
| +7      | 0111                | 0111               | 0111               | 1111           | 1110           |
| +6      | 0110                | 0110               | 0110               | 1110           | 1101           |
| +5      | 0101                | 0101               | 0101               | 1101           | 1100           |
| +4      | 0100                | 0100               | 0100               | 1100           | 1011           |
| +3      | 0011                | 0011               | 0011               | 1011           | 1010           |
| +2      | 0010                | 0010               | 0010               | 1010           | 1001           |
| +1      | 0001                | 0001               | 0001               | 1001           | 1000           |
| +0      | 0000                | 0000               | 0000               | 1000           | 0111           |
| -0      | 1000                | 1111               | 0000               | -              | -              |
| -1      | 1001                | 1110               | 1111               | 0111           | 0110           |
| -2      | 1010                | 1101               | 1110               | 0110           | 0101           |
| -3      | 1011                | 1100               | 1101               | 0101           | 0100           |
| -4      | 1100                | 1011               | 1100               | 0100           | 0011           |
| -5      | 1101                | 1010               | 1011               | 0011           | 0010           |
| -6      | 1110                | 1001               | 1010               | 0010           | 0001           |
| -7      | 1111                | 1000               | 1001               | 0001           | 0000           |
| -8      | -                   | -                  | 1000               | 0000           | -              |

STUDENTS-HUB.com

#### Uploaded By: anonymous

11

### Sign Bit

- Highest bit indicates the sign
- 1 = negative
- $\bullet$  0 = positive



- For Hexadecimal Numbers, check most significant digit
  - If highest digit is > 7, then value is negative
- Examples: 8A and C5 are negative bytes
- B1C42A00 is a negative word (32-bit signed integer)
- Problems
  - Need to consider both sign and magnitude in arithmetic
  - Two representations of zero (+0 and -0) in sign-mag & 1's comp

### **Biased Representation**

- Other type of binary number representations
- A fixed value called Bias is added for the binary value
- Typically, the bias equals (2<sup>k-1</sup>-1), where K is the number of bits in the binary number.
- e.g for 4 bit representation,
  - The bias value= 2<sup>4-1</sup>-1= 7
  - Representation of +8 => 8+7=15 => 1111
  - Representation of -7 => -7+7=0 => 0000

# **Two's Complement Representation**

- Positive numbers
  - Signed value = Unsigned value
- Negative numbers
  - Signed value = Unsigned value 2<sup>n</sup>
  - n = number of bits
- Negative weight for MSB
- Another way to obtain the signed value is to assign a negative weight to most-significant bit



| 8-bit Binary value | Unsigned<br>value | Signed value |
|--------------------|-------------------|--------------|
| 0000000            | 0                 | 0            |
| 0000001            | 1                 | +1           |
| 00000010           | 2                 | +2           |
|                    |                   |              |
| 01111110           | 126               | +126         |
| 01111111           | 127               | +127         |
| 10000000           | 128               | -128         |
| 10000001           | 129               | -127         |
|                    |                   |              |
| 11111110           | 254               | -2           |
| 11111111           | 255               | -1           |

### Forming the Two's Complement

| starting value                           | 00100100 = +36 |
|--|----------------|
| step1: reverse the bits (1's complement) | 11011011       |
| step 2: add 1 to the value from step 1   | + 1            |
| sum = 2's complement representation      | 11011100 = -36 |

Sum of an integer and its 2's complement must be zero:  $00100100 + 11011100 = 00000000 (8-bit sum) \Rightarrow$  Ignore Carry

Another way to obtain the 2's complement:Binary ValueStart at the least significant 1= 00100100 significant 1Leave all the 0s to its right unchanged2's ComplementComplement all the bits to its left= 11011100

### Sign Extension

- **Step 1**: Move the number into the lower-significant bits
- **Step 2**: Fill all the remaining higher bits with the sign bit
- This will ensure that both magnitude and sign are correct

#### Examples

◇ Sign-Extend 10110011 to 16 bits
10110011 = -77 → 1111111111100110011 = -77
◇ Sign-Extend 01100010 to 16 bits
01100010 = +98 → 00000000001100010 = +98
Infinite 0s can be added to the left of a positive number
Infinite 1s can be added to the left of a negative number

# **Ranges of Signed Integers**

For *n*-bit signed integers: Range is  $-2^{n-1}$  to  $(2^{n-1} - 1)$ 

Positive range: 0 to  $2^{n-1} - 1$ 

Negative range:  $-2^{n-1}$  to -1

| Storage Type | Signed Range                     | Powers of 2                               |  |
|--------------|----------------------------------|---|--|
| Byte         | -128 to +127                     | -2 <sup>7</sup> to (2 <sup>7</sup> - 1)   |  |
| Half Word    | -32,768 to +32,767               | -2 <sup>15</sup> to (2 <sup>15</sup> - 1) |  |
| Word         | -2,147,483,648 to +2,147,483,647 | -2 <sup>31</sup> to (2 <sup>31</sup> - 1) |  |
| Double Word  | -9,223,372,036,854,775,808 to    | $263 \pm (263 \pm 1)$                     |  |
| Double word  | +9,223,372,036,854,775,807       | $-2^{00}$ to $(2^{00} - 1)$               |  |

Practice: What is the range of signed values that may be stored in 20 bits?



**Figure 10.5 Geometric Depiction of Twos Complement Integers** 

STUDENTS-HUB.com

### Negation

Twos complement operation

- Take the Boolean complement of each bit of the integer (including the sign bit)
- Treating the result as an unsigned binary integer, add 1

+18 = 00010010 (twos complement) bitwise complement = 11101101  $\frac{+ 1}{11101110} = -18$ 

The negative of the negative of that number is itself:

```
-18 = 11101110 \text{ (twos complement)}
bitwise complement = 00010001
+ 1
```

00010010 = +18

Uploaded By: anonymous

### **Negation Special Case 1**

0 = 0000000 (twos complement)Bitwise complement = 11111111 Add 1 to LSB + 1Result 10000000

Overflow is ignored, so:

-0 = 0

### **Negation Special Case 2**

-128 = 10000000 (twos complement)

Bitwise complement = 01111111

Add 1 to LSB

Result

10000000

So:

-(-128) = -128 X

Monitor MSB (sign bit)

It should change during negation

#### Table 10.1 Characteristics of Twos Complement Representation and Arithmetic

| Range                                | $-2^{n-1}$ through $2^{n-1} - 1$  |
|--------------------------------------|---|
| Number of Representations<br>of Zero | One   |
| Negation                             | Take the Boolean complement of each bit of the corresponding positive number, then add 1 to the resulting bit pattern viewed as an unsigned integer.        |
| Expansion of Bit Length              | Add additional bit positions to the left and fill in with the value of the original sign bit.   |
| Overflow Rule                        | If two numbers with the same sign (both positive or both nega-<br>tive) are added, then overflow occurs if and only if the result has<br>the opposite sign. |
| Subtraction Rule                     | To subtract $B$ from $A$ , take the twos complement of $B$ and add it to $A$ .  |

# **Character Storage**

#### Character sets

- Standard ASCII: 7-bit character codes (0 127)
- Extended ASCII: 8-bit character codes (0 255)
- Unicode: 16-bit character codes (0 65,535)
- Unicode standard represents a universal character set
  - Defines codes for characters used in all major languages
  - Used in Windows-XP: each character is encoded as 16 bits
- UTF-8: variable-length encoding used in HTML
  - Encodes all Unicode characters
  - Uses 1 byte for ASCII, but multiple bytes for other characters
- Null-terminated String
  - Array of characters followed by a NULL character

# **Binary Addition**

- Start with the least significant bit (rightmost bit)
- Add each pair of bits
- Include the carry in the addition, if present



### **Binary Subtraction**

■ When subtracting A – B, convert B to its 2's complement

Add A to (-B)



Final carry is ignored, because

- Negative number is sign-extended with 1's
- You can imagine infinite 1's to the left of a negative number
- Adding the carry to the extended 1's produces extended zeros

STUDENTS-HUB.com

| 1001 = -7 + 0101 = 5 = -2<br>1110 = -2<br>(a) (-7) + (+5) | 1100 = -4 + 0100 = 4 10000 = 0 (b) (-4) + (+4)                   |
|---|--|
| 0011 = 3 + 0100 = 4 0111 = 7 (c) (+3) + (+4)              | 1100 = -4 + 1111 = -1 11011 = -5 (d) (-4) + (-1)                 |
| $0101 = 5 + 0100 = 4 \\ 1001 = Overflow$ (e) (+5) + (+4)  | 1001 = -7 + 1010 = -6 = -6<br>10011 = Overflow<br>(f)(-7) + (-6) |

**Figure 10.3** Addition of Numbers in Twos Complement Representation

**OVERFLOW RULE:** If two numbers are added, and they are both positive or both negative, then overflow occurs if and only if the result has the opposite sign.

| $\begin{array}{rcrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$                                | $\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$       |
|--|--|
| (a) $M = 2 = 0010$   | (b) $M = 5 = 0101$   |
| S = 7 = 0111   | S = 2 = 0010   |
| -S = 1001  | -S = 1110  |
| $1011 = -5 + \frac{1110}{1001} = -2 -7$  | $\begin{array}{rcrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$      |
| (c) $M = -5 = 1011$  | (d) $M = 5 = 0101$   |
| S = 2 = 0010   | S = -2 = 1110  |
| -S = 1110  | -S = 0010  |
| $ \begin{array}{rcl} 0111 &=& 7 \\ + & 0111 &=& 7 \\ 1110 &=& Overflow \end{array} $ | 1010 = -6<br>+ <u>1100</u> = -4<br><u>10110</u> = Overflow |
| (e) $M = 7 = 0111$   | (f) $M = -6 = 1010$  |
| S = -7 = 1001  | S = 4 = 0100   |
| -S = 0111  | -S = 1100  |

**Figure 10.4** Subtraction of Numbers in Twos Complement Representation (M - S)

**SUBTRACTION RULE:** To subtract one number (subtrahend) from another (minuend), take the twos complement (negation) of the subtrahend and add it to the minuend.

STUDENTS-HUB.com



**SW** = Switch (select addition or subtraction)

Figure 10.6 Block Diagram of Hardware for Addition and STUDENTS-HUB.com Upload 28

# **Unsigned Multiplication**

Paper and Pencil Example:

| Multiplicand<br>Multiplier | 1100 <sub>2</sub><br>× 1101 <sub>2</sub> | = 12<br>= 13   |
|----------------------------|--|--|
|                            | 1100<br>0000<br>1100<br>1100             | Binary multiplication is easy<br>0 × multiplicand = 0<br>1 × multiplicand = multiplicand |
| Product                    | 10011100.                                | = 156  |

m-bit multiplicand × n-bit multiplier = (m+n)-bit product

Accomplished via shifting and addition

Consumes more time and more chip area



Figure 10.7 Multiplication of Unsigned Binacket Bygersonymous





Figure 10.8 Hardware Implementation of Unsigned Binary Multiplication

#### STUDENTS-HUB.com

# **Multiplying Negative Numbers**

This does not work!

#### Solution 1

- Convert to positive if required
- Multiply as above
- If signs were different, negate answer

#### Solution 2

Booth's algorithm



#### Booth's algorithm

https://www.grahn.us/projects/booths-algorithm/

|      | -    | -        |      |                  |          |                |
|------|------|----------|------|------------------|----------|----------------|
| A    | Q    | $Q_{-1}$ | M    |                  |          |                |
| 0000 | 0011 | 0        | 0111 | Initial value    | <b>S</b> |                |
| 1001 | 0011 | 0        | 0111 | A←A – M          | 2        | First          |
| 1100 | 1001 | 1        | 0111 | Shift            | 5        | cycle          |
| 1110 | 0100 | 1        | 0111 | Shift            | }        | Second cycle   |
| 0101 | 0100 | 1        | 0111 | A←A + M<br>Shift | }        | Third<br>cvcle |
| 0010 | 1010 | 0        | 0111 | omit             | <u>ر</u> | Fourth         |
| 0001 | 0101 | 0        | 0111 | Shift            | Ş        | cycle          |

Figure 10.13 Example of Booth's Algarithmy (Znon3)mous

| 3  |                                   |   |   |
|--|-----------------------------------|---|---|
| 0111<br><u>´0011</u><br>11111001<br>0000000<br>000111<br>00010101        | (0)<br>1-0<br>1-1<br>0-1<br>(21)  | 0111<br><u>´1101</u><br>11111001<br>0000111<br>111001<br>11101011 | (0)<br>1-0<br>0-1<br>1-0<br>(-21)   |
| (a) $(7)$  | (2) (21)                          | (1-) (7)  | (2) $(21)$  |
| (a) (7)  | (3) = (21)                        | (D)(/)  | (-3) = (-21)  |
| ALCOME THE OWNER DURING THE PARTY AND A 1272                             |                                   | CONTRACTOR OF CONTRACTOR OF CONTRACTOR                            | strate of the second |
| 1001<br><u>´0011</u><br>00000111<br>0000000<br><u>111001</u><br>11101011 | (0)<br>1-0<br>1-1<br>0-1<br>(-21) | 1001<br><u>´1101</u><br>00000111<br>1111001<br>000111<br>00010101 | (0)<br>1-0<br>0-1<br>1-0<br>(21)  |

(c) (-7) (3) = (-21)

(d) (-7) (-3) = (21)

#### Figure 10.14 Examples Using Booth's Algorithm

STUDENTS-HUB.com

### How it works

Consider a positive multiplier consisting of a block of 1s surrounded by 0s. For example, 00111110. The product is given by :

$$M \times "0\ 0\ 1\ 1\ 1\ 1\ 1\ 0\ " = M \times (2^5 + 2^4 + 2^3 + 2^2 + 2^1) = M \times 62$$

• where M is the multiplicand.

The number of operations can be reduced to two by rewriting the same as

$$M \times "0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ -1\ 0" = M \times (2^6 - 2^1) = M \times 62$$

• Note that:

$$2^{n} + 2^{n-1} + ... + 2^{n-k} = 2^{n+1} - 2^{n-k}$$

### How it works

- So, the product can be generated by one addition and one subtraction
- In Booth's algorithm
  - perform subtraction when the first 1 of the block is encountered (1 - 0)
  - perform addition when the last 1 of the block is encountered (0 1)
- (1 0) and (0 1) are observed from Q<sub>0</sub> Q<sub>-1</sub> (see previous example)

# Division

- More complex than multiplication
- Negative numbers are really bad!
- Based on long division

36



**Figure 10.15 Example of Division of Unsigned Binary Integers** 

STUDENTS-HUB.com



### **Real Numbers**

- Numbers with fractions
- Could be done in pure binary
   1001.1010 = 2<sup>3</sup> + 2<sup>0</sup> + 2<sup>-1</sup> + 2<sup>-3</sup> = 9.625
- Where is the binary point?
- Fixed?
  - Very limited
- Moving?
  - How do you show where it is?

39

### **Exponential Notation**

The following are equivalent representations of 1,234

| 123,400 | • 0   | Х | 10-2            |
|---------|-------|---|-----------------|
| 12,340  | . 0   | Х | 10-1            |
| 1,234   | .0    | Х | 100             |
| 123     | . 4   | Х | 101             |
| 12      | .34   | Х | 10 <sup>2</sup> |
| 1       | .234  | Х | 10 <sup>3</sup> |
| 0       | .1234 | Х | 104             |
|         |       |   |                 |

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

# **Floating Point**

- An IEEE Std 754 floating point representation consists of
  - A Sign Bit (no surprise)
  - An Exponent ("times 2 to the what?")
  - Mantissa ("Significand"), which is assumed to be 1.xxxxx (thus, one bit of the mantissa is implied as 1)
  - This is called a normalized representation
- So a mantissa = 0 really is interpreted to be 1.0, and a mantissa of all 1111 is interpreted to be 1.1111
- Special cases are used to represent denormalized mantissas (true mantissa = 0), NaN, etc.



#### IEEE 754 Floating-Point Standard

| TYPES                                | SIGN        | BIASED<br>EXPONENT | NORMALISED<br>MANTISA | BIAS                       |
|--------------------------------------|-------------|--------------------|-----------------------|----------------------------|
| Single precision                     | l(31st bit) | 8(30-23)           | 23(22-0)              | 127                        |
| Double precision<br>STUDENTS-HUB.con | 1(63rd bit) | 11(62-52)          | 52(51-0)<br>Uplo      | 1023<br>aded By: anonymous |



85.125
85 = 1010101
0.125 = 001
85.125 = 1010101.001
=1.010101001 × 2^6
sign = 0

1. Single precision: biased exponent 127+6=133 133 = 10000101 Normalised mantisa = 010101001 we will add 0's to complete the 23 bits

The IEEE 754 Single precision is: = 0 10000101 01010010000000000000000 This can be written in hexadecimal form **42AA4000** 

#### STUDENTS-HUB.com

What number is represented by the single precision float 11000000101000...00

1100000101000...00

-S = 1

- Fraction = 01000...00<sub>2</sub>
- Fxponent = 10000001<sub>2</sub> = 129
- $x = (-1)^1 \times (1 + 01_2) \times 2^{(129 127)}$
- $= (-1) \times 1.25 \times 2^{2}$

Represent -0.75

- $--0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
- -S = 1
- Fraction = 1000...00<sub>2</sub>
- Exponent = -1 + Bias
  - Single: -1 + 127 = 126 = 01111110<sub>2</sub>
  - Double: -1 + 1023 = 1022 = 01111111102

Single: 1011111101000...00 Double: 101111111101000...00 45

### **Converting from Floating Point**

E.g., What decimal value is represented by the following 32bit floating point number?

**C17B0000**<sub>16</sub>



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

### **Converting from Floating Point**

- E.g., What decimal value is represented by the following 32bit floating point number?
  - **C17B0000**<sub>16</sub>



-1111.1011,



### **Converting to Floating Point**

E.g., Express 36.5625<sub>10</sub> as a 32-bit floating point number (in hexadecimal)



Uploaded By: anonymous

### **Converting to Floating Point**

Step 5

E.g., Express 36.5625<sub>10</sub> as a 32-bit floating point number (in hexadecimal)



#### Figure 10.21 IEEE 754 Formats

#### STUDENTS-HUB.com

Uploaded By: anonymous

50

# FP Arithmetic +/-

- Check for zeros
- Align significands (adjusting exponents)
- Add or subtract significands
- Normalize result



Figure 10.22 Floating-Point Addition and Subtraction (Z-X ± Y)

STUDENTS-HUB.com



Perform 0.5 + (-0.4375)

 $0.5 = 0.1 \times 2^{0} = 1.000 \times 2^{-1}$  (normalised)

-0.4375 = -0.0111 × 2<sup>0</sup> = -1.110 × 2<sup>-2</sup> (normalised)

1. Rewrite the smaller number such that its exponent matches with the exponent of the larger number.

 $-1.110 \times 2^{-2} = -0.1110 \times 2^{-1}$ 

2. Add the mantissas:

 $1.000 \times 2^{-1} + -0.1110 \times 2^{-1} = 0.001 \times 2^{-1}$ 

0.001 × 2<sup>-1</sup> = 1.000 × 2<sup>-4</sup> -126 <= -4 <= 127 ===> No overflow or underflow

The sum fits in 4 bits so rounding is not required

Check: 1.000 × 2<sup>-4</sup> = 0.0625 which is equal to 0.5 - 0.4375

3. Normalise the sum, checking for overflow/underflow:

#### STUDENTS-HUB.com

4. Round the sum:

#### **Step 1: Decompose Operands**

Add the floating point numbers 3.75 and 5.125 to get 8.875 by directly manipulating the numbers in IEEE format.

|                        | Sign | Exponent                          | Mantissa                                |
|------------------------|------|-----------------------------------|---|
| Value:                 | +1   | 21                                | 1.875                                   |
| Encoded as:            | 0    | 128                               | 7340032                                 |
| Binary:                |      |                                   |   |
|                        |      | Decimal Representation            | 3.75                                    |
|                        |      | Binary Representation             | 010000000111000000000000000000000000000 |
|                        |      | Hexadecimal Representation        | 0x40700000                              |
|                        |      | After casting to double precision | 3.75                                    |
|                        |      |                                   |   |
|                        | Sign | Exponent                          | Mantissa                                |
| Value:                 | +1   | 22                                | 1.28125                                 |
| Encoded as:<br>Binary: | 0    |                                   |   |
|                        |      | Decimal Representation            | 5.125                                   |
|                        |      | Binary Representation             | 010000001010010000000000000000000000000 |
|                        |      | Hexadecimal Representation        | 0x40a40000                              |
|                        |      | After casting to double precision | 5.125                                   |

- For 3.75, the sign bit is 0 (+), the exponent is 128 (1 unbiased), the mantissa (including the implicit 1 shown in bold) is: 0000 0000 1111 00000 0000 0000 0000 = 0x00f00000
- For 5.125, the sign bit is 0 (+), the exponent is 129 (2 unbiased), the mantissa (including the implicit 1 shown in bold) is: 0000 0000 1010 0100 0000 0000 0000 = 0x00a40000
   STUDENTS-HUB.com

#### **Step 2: Equalizing Operand Exponents**

- If the first exponent is smaller than the second, we shift the first mantissa to the right and add the absolute difference in exponents to the first exponent. If vice versa, we do the same to the second mantissa and exponent.
- For this example the first exponent is 128, second exponent is 129, absolute difference is 1, so first exponent is smaller, so we must adjust the first mantissa and exponent, and leave the second mantissa and exponent unchanged.
  - Shift first mantissa right by 1: 0x00f00000 >> 1 = 0x00780000
  - Increase the first exponent by 1: 128 + 1 = 129

#### <u>Step 3: Convert operands from signed magnitude to 2's</u> <u>complement</u>

For each operand that is negative, convert the mantissa to 2's complement by inverting the bits and adding 1. Neither operand is negative in this example, so nothing needs to be done.

#### **Step 4: Add Mantissas**

- Both operands have an exponent of 129, so we can just add mantissas to get a positive result with the same exponent.

#### <u>Step 5: Convert result from 2's complement to signed</u> <u>magnitude</u>

If the result is negative, convert the mantissa back to signed magnitude by inverting the bits and adding 1. The result is positive in this example, so nothing needs to be done.

STUDENTS-HUB.com

#### **Step 6: Normalize Result**

- Because the leftmost 1 bit is not in the right place, we must shift the mantissa right or left to put it back into the IEEE format, and adjust the exponent accordingly.
- If the leftmost 1 bit is left of bit 23, we must shift the mantissa to the right and increase the exponent. If the leftmost 1 bit is at bit 23, there is no normalization required. If the leftmost 1 bit is right of bit 23, we must shift the mantissa to the left and decrease the exponent. In the example, we see the first case:

#### BEFORE NORMALIZATION

```
Sign of result = 0
Exponent of result = 129
Mantissa of result = 0000 0001 0001 1100 0000 0000 0000 = 0x011c0000
```

The leftmost 1 bit is bit 24, so we must shift the mantissa right by 1 and add 1 to the exponent:

#### AFTER NORMALIZATION

#### Step 7: Compose Result

| Value        | Sign       | Exponent                          | Mantissa                                |
|--------------|------------|-----------------------------------|---|
| value:       | +1         | 25                                | 1.109375                                |
| Encoded as:  | 0          | 130                               | 917504                                  |
| Binary:      |            |                                   |   |
|              |            | Decimal Representation            | 8.875                                   |
|              |            | Binary Representation             | 010000010000111000000000000000000000000 |
|              |            | Hexadecimal Representation        | 0x410e0000                              |
|              |            | After casting to double precision | 8.875                                   |
| To summarize | the result | of the addition:                  |   |

3.75 (0x40700000) = +5.125 (0x40a40000) = 8.875 (0x410e0000)

Activate Windo Go to Settings to ac

# FP Arithmetic x/÷

- Check for zero
- Add/subtract exponents
- Multiply/divide significands (watch sign)
- Normalize
- Round
- All intermediate results should be in double length storage

59



Figure 10.23 Floating-Point Multiplication (Z← X× Y)

STUDENTS-HUB.com

Uploaded By: anonymous

60

 $1.000 \times 2^{-1} \times -1.110 \times 2^{-2}$ 

1. Add the biased exponents

(-1 + 127) + (-2 + 127) - 127 = 124 ===> (-3 + 127)

2. Multiply the mantissas

3. Normalise (already normalised)

At this step check for overflow/underflow by making sure that

-126 <= Exponent <= 127

1 <= Biased Exponent <= 254

4. Round the result (no change)

5. Adjust the sign.

Since the original signs are different, the result will be negative



Figure 10.24 Floating-Point Division (Z← X/Y)

STUDENTS-HUB.com