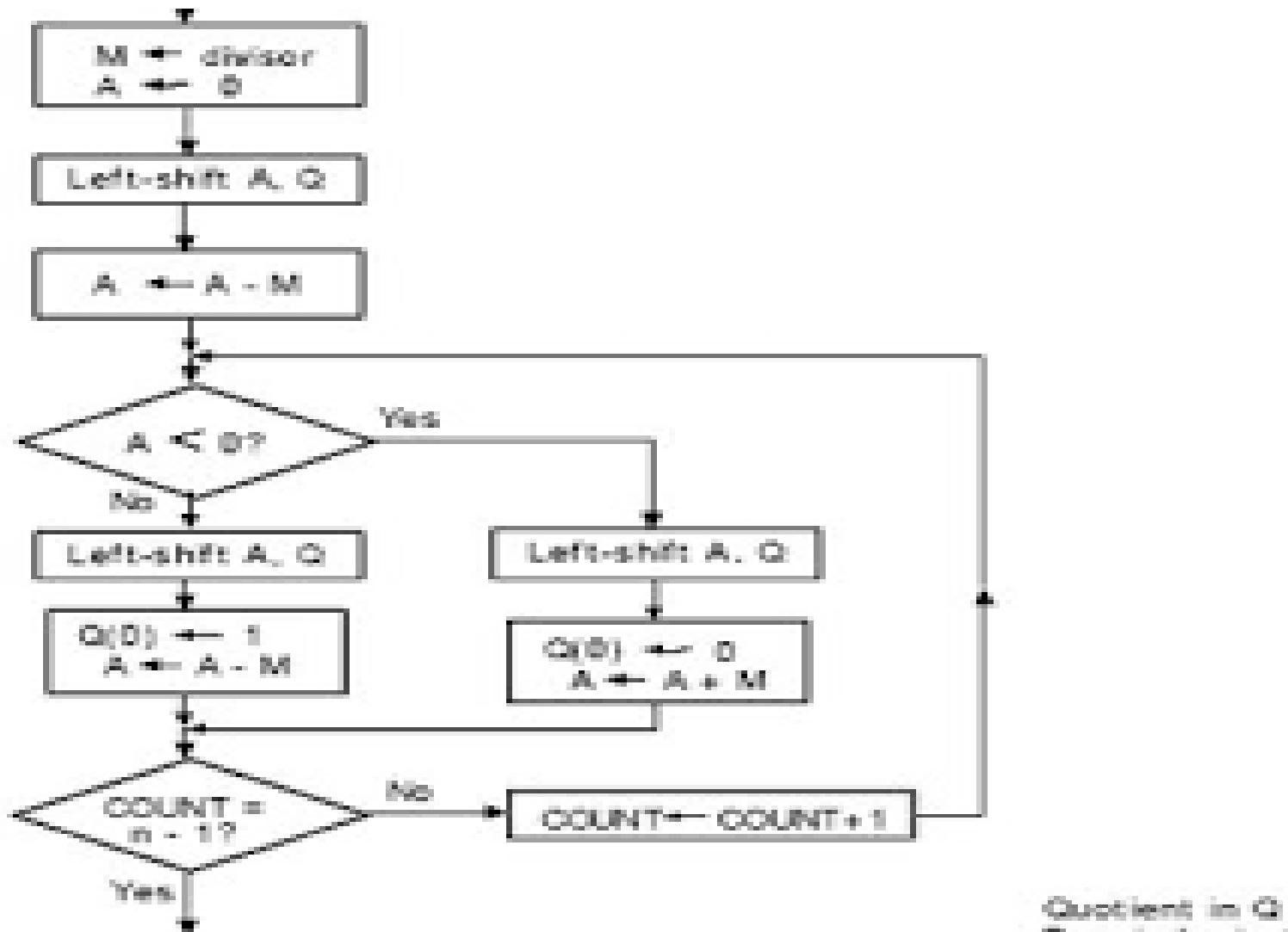
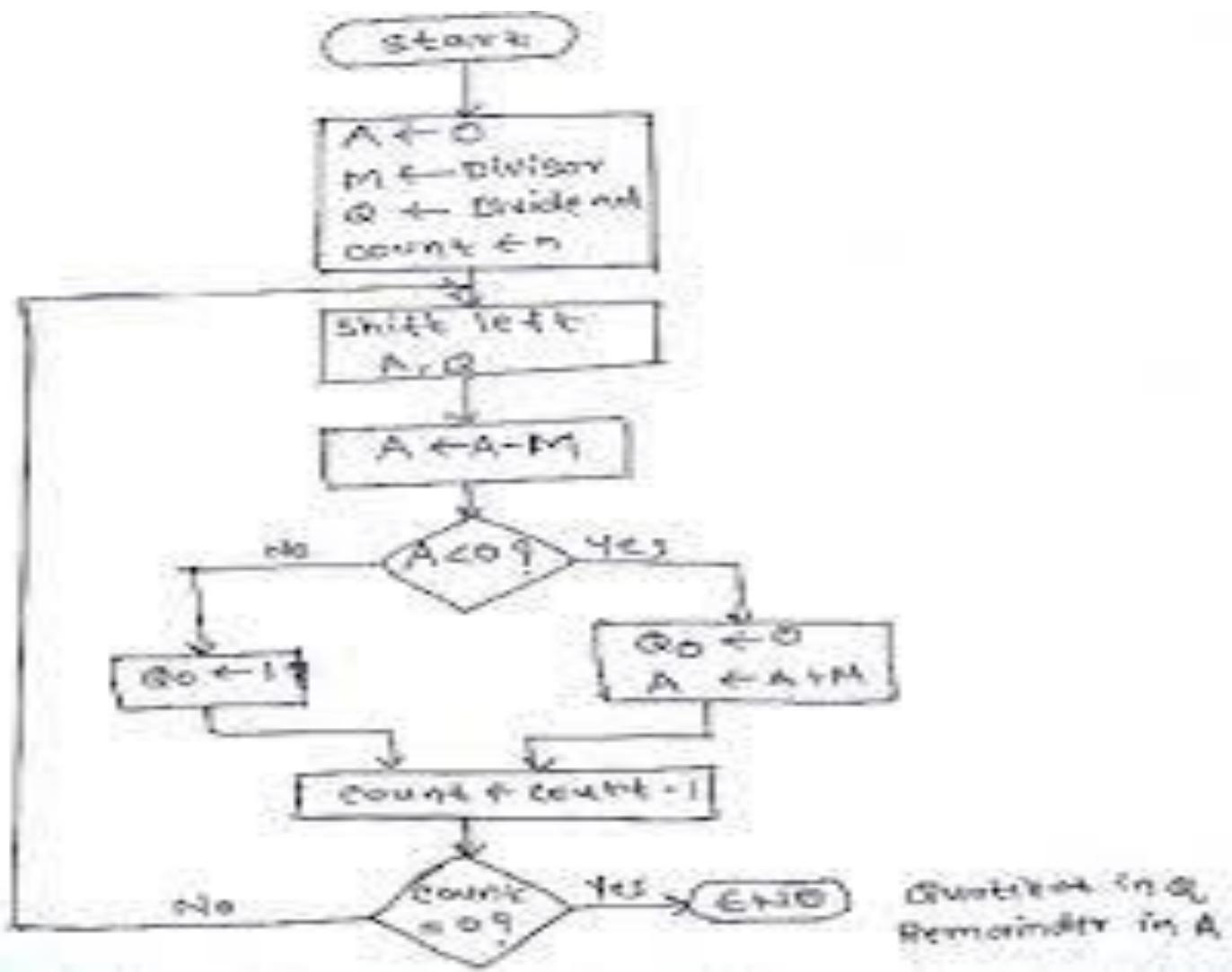


Booths division –restoring algorithm

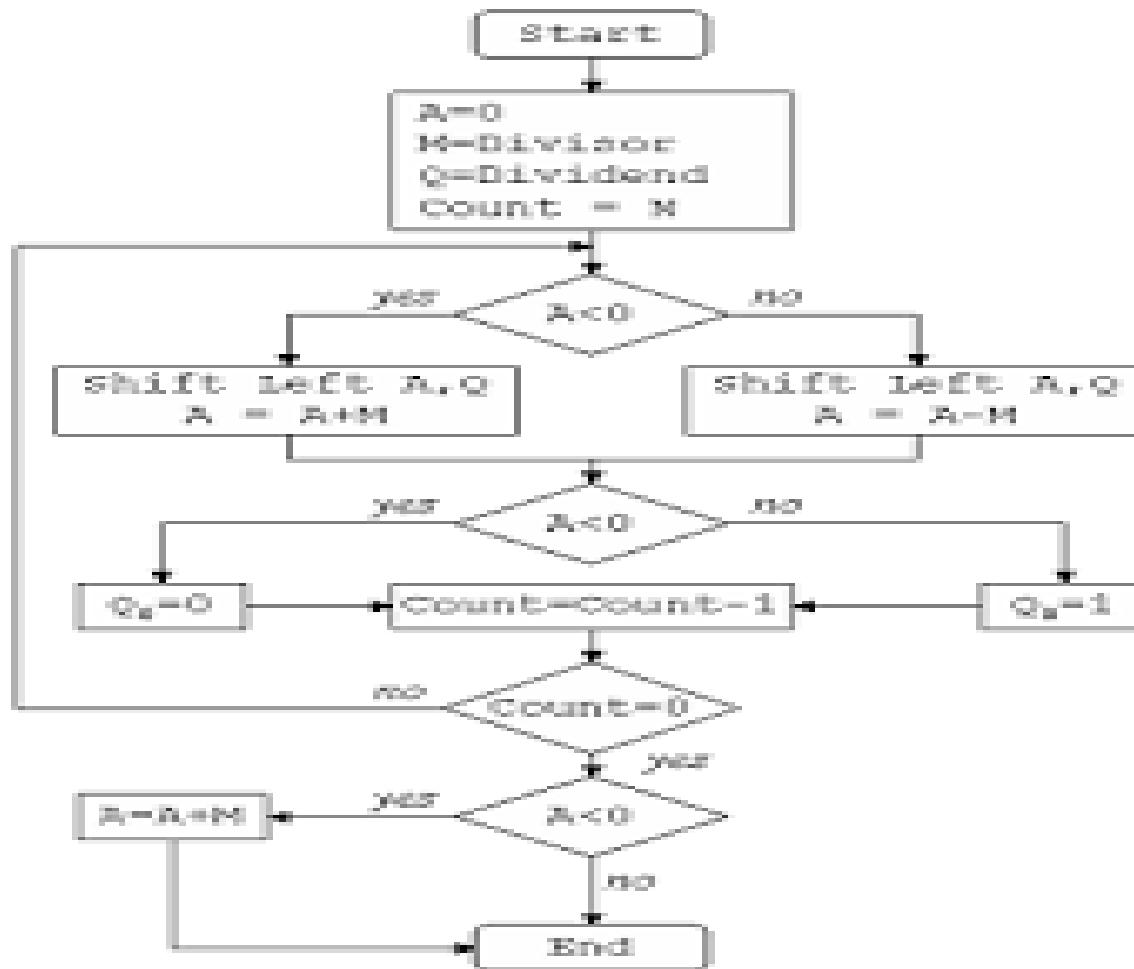


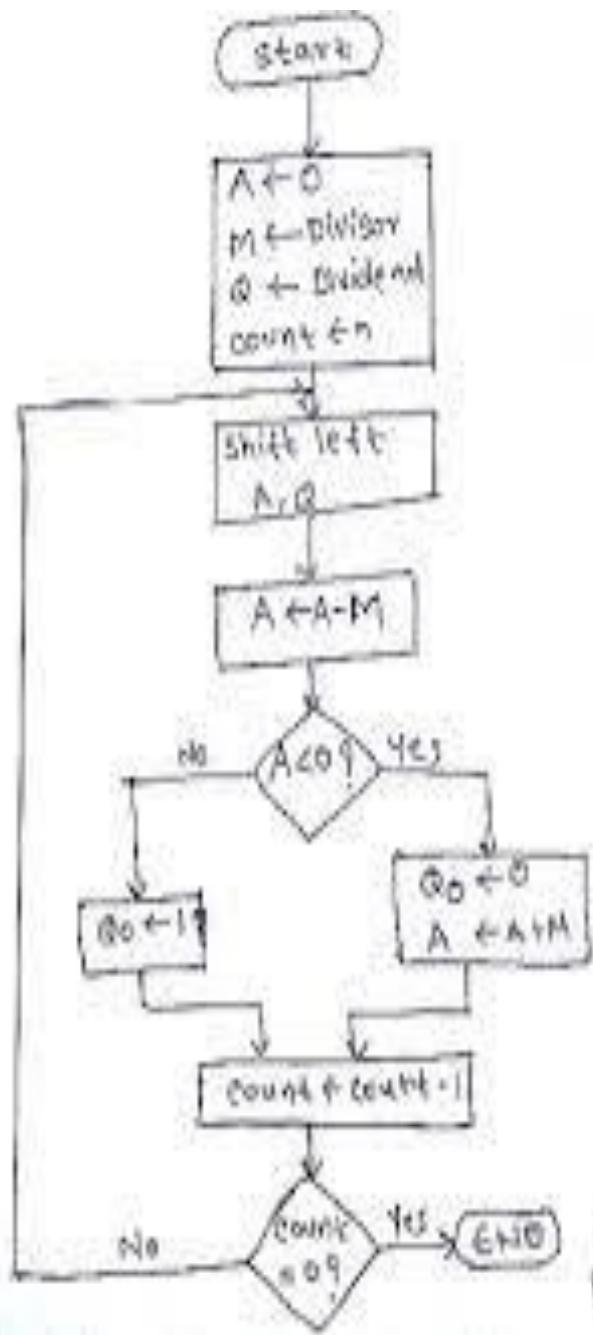
Booth's division –restoring algorithm



Quotient in Q
Remainder in A

Booths division –non-restoring algorithm





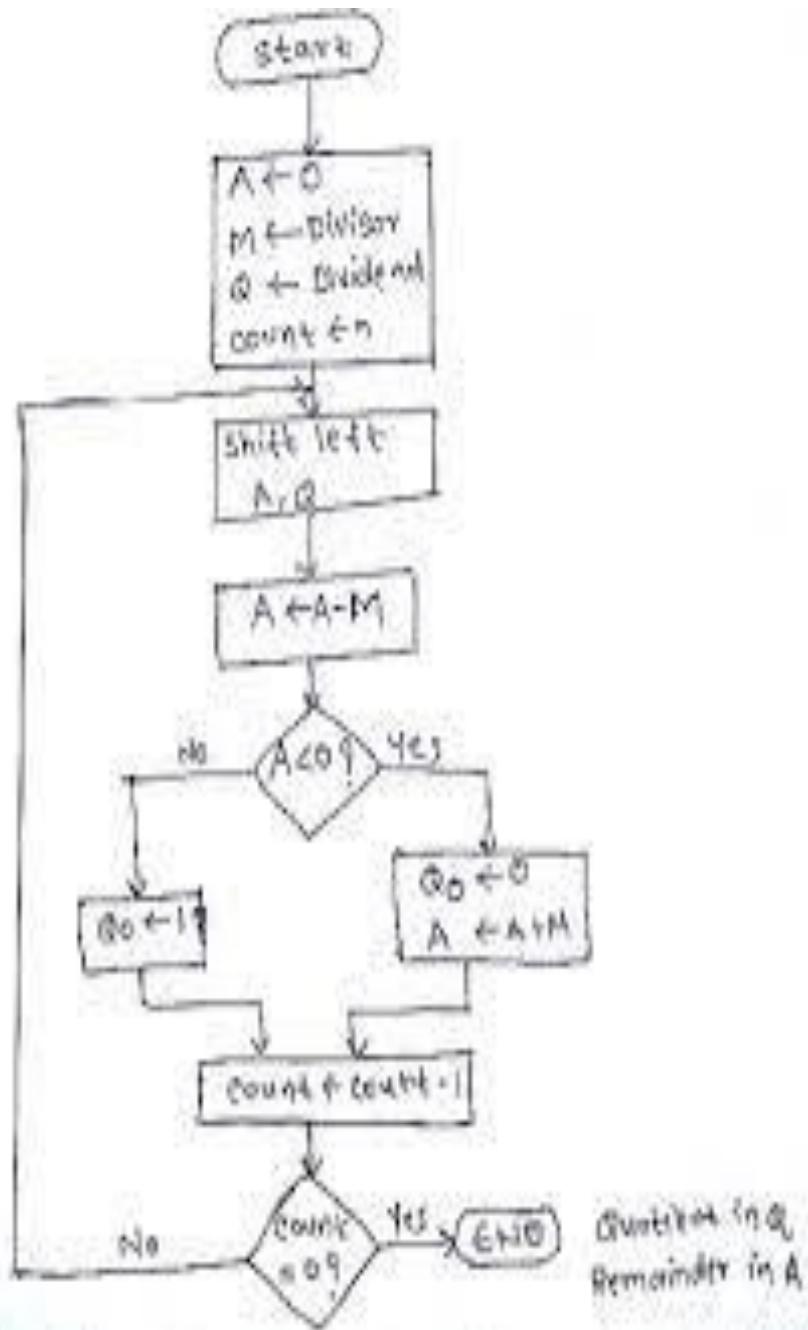
	10
10	1000
	11
	10
Initially	0 0 0 0 0 1 0 0 0
	0 0 0 1 1
Subtract	0 0 0 0 1 0 0 0
	1 1 1 0 1
Set q0	1 1 1 1 0
Restore	1 1 0 0 0 0 0
	0 0 0 0 1 0 0 0 0 0
Shift	0 0 0 1 0 0 0 0
Subtract	1 1 1 0 1
Set q1	1 1 1 1 1
Restore	1 1 0 0 0 0 0
	0 0 0 1 0 0 0 0 0 0
Shift	0 0 1 0 0 0 0 0
Subtract	1 1 0 1
Set q2	0 0 0 0 1
Restore	0 0 0 1 0 0 0 0 0 0
	0 0 1 0 0 0 0 0 0 0
Shift	0 0 0 1 0 0 0 0 0 1
Subtract	1 1 1 0 1 0 0 1
Set q3	1 1 1 1 1
Restore	1 1 0 0 1 0
	0 0 0 1 0 0 0 1 0
	Remainder
	Quotient

First cycle

Second cycle

Third cycle

Fourth cycle



10

Initially 0 0 0 0 0 1 0 0 0

Был в в в в в | в в в

Subtract 1 1 0 1

See Fig. 1

Restore

1 / 1

Shift 0 0 0 1 0 0 0 0

Subtract: 1 1 1 0 1

Zentrierte Φ_0  | | | | |

Restore

0 0 0 1 0 0 0 0

Shift 0 0 1 0 0 0 0 0

$$\begin{array}{r} \text{Subtract} \\ - 12 \\ \hline \end{array}$$

Set η_0 to 0 0 0 0 1

10.1002/anie.201907002

Software 1 1 1 2 1 2 2 1

Review

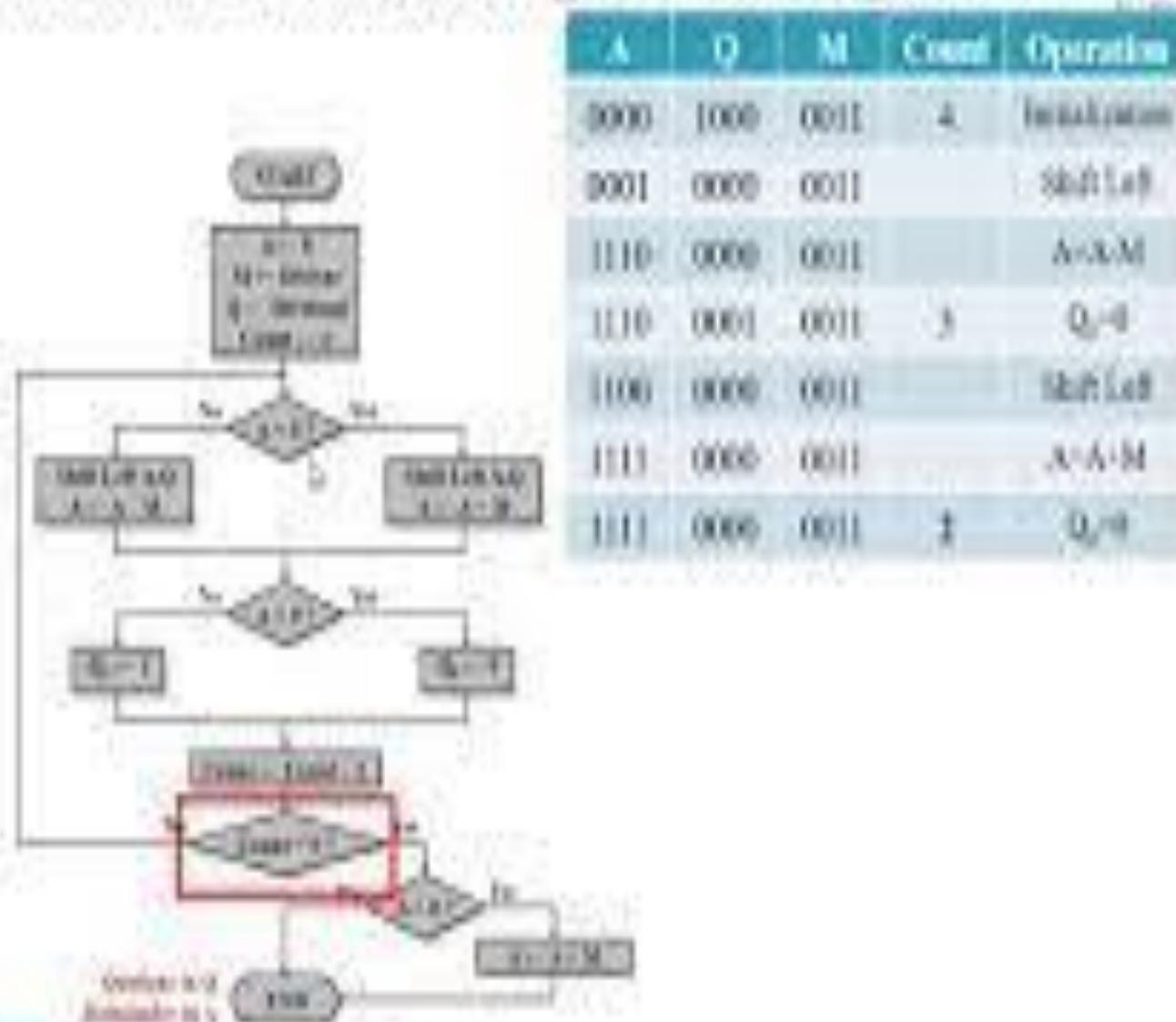
For more information about the study, please contact Dr. Michael J. Hwang at (310) 206-6500 or via email at mhwang@ucla.edu.

... [View Details](#) [Edit](#) [Delete](#) [Print](#)

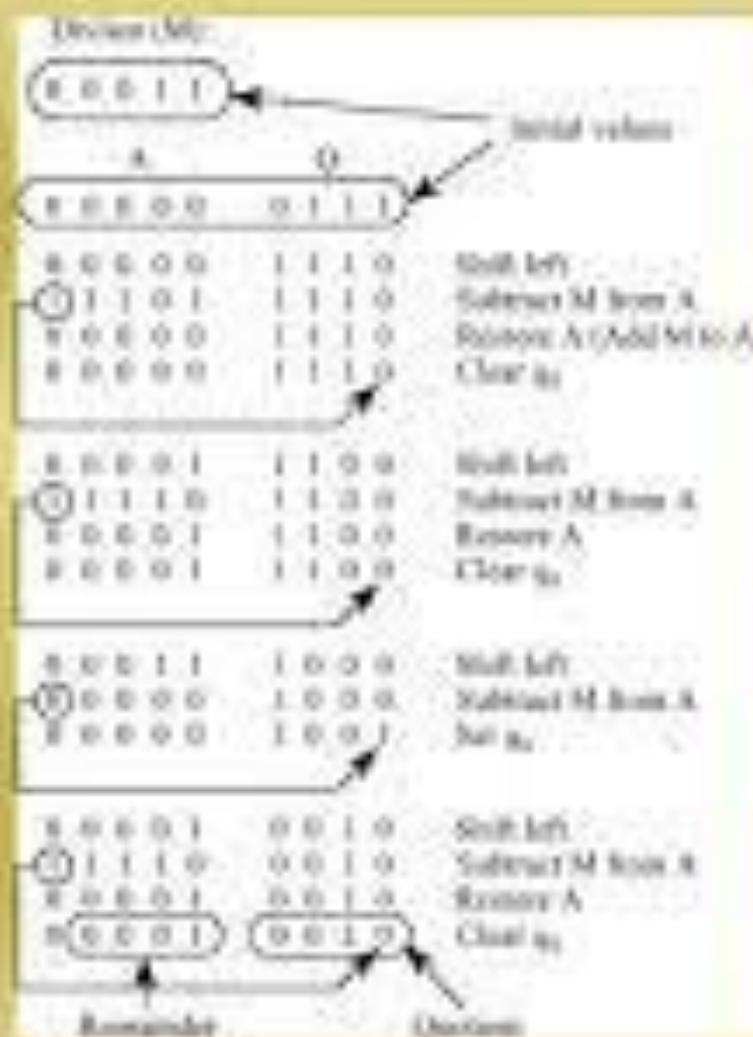
1

1

Booth's Non-Restoring Division Algorithm Deo's



Division Example Using Serial Divider

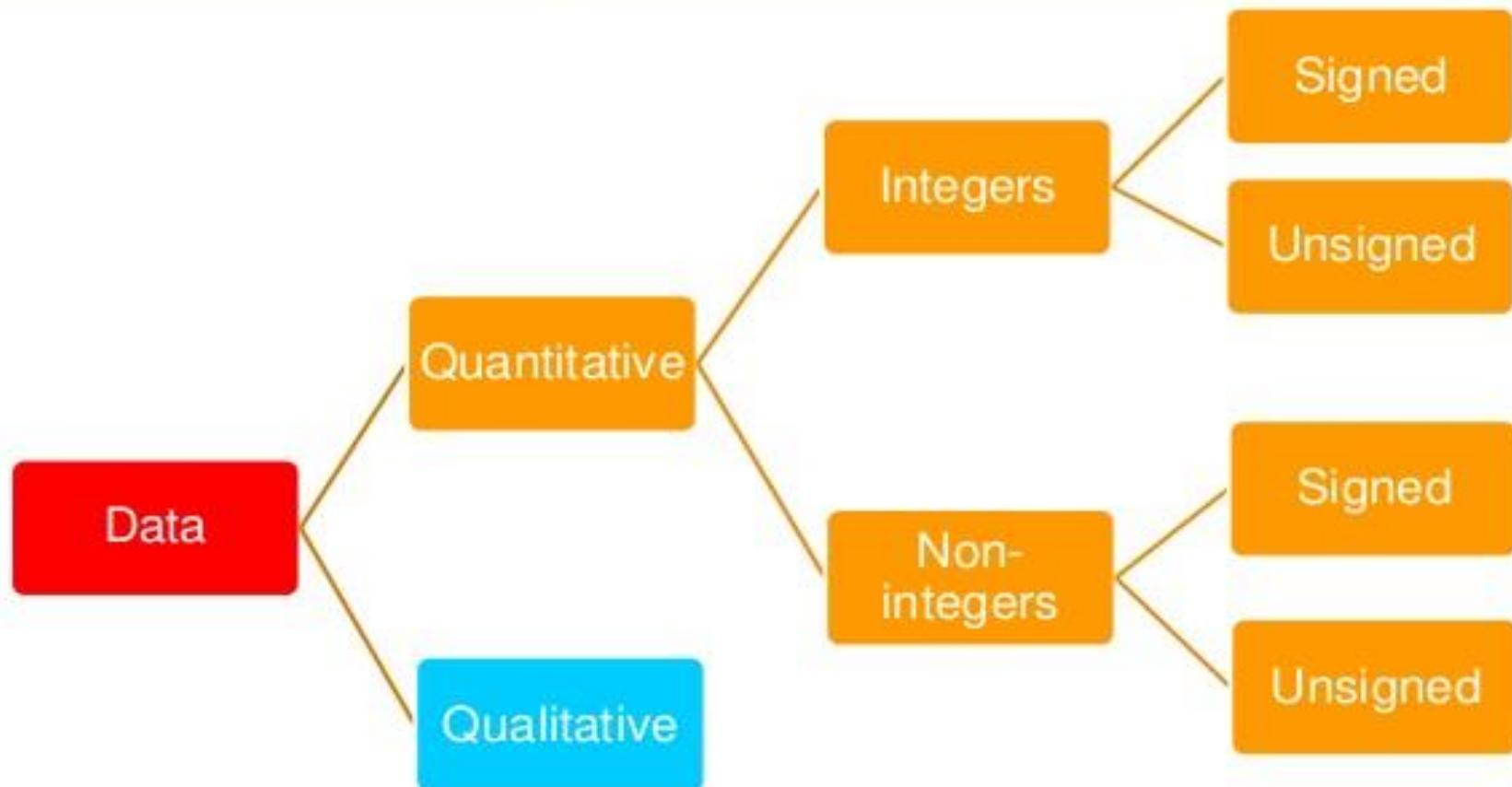


Data Representation

Outline

- Representing numbers
 - Unsigned
 - Signed
 - Floating point
- Representing characters & symbols
 - ASCII
 - Unicode

Data Representation (Cont.)



Number Systems

- Decimal number system
 - 0, 1, 2, 3, 4, 5, 6, 7, 8, 9
- Binary number system
 - 0, 1
- Octal number system
 - 0, 1, 2, 3, 4, 5, 6, 7
- Hexadecimal number system
 - 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F

Quantitative Numbers

- Integers

- Unsigned 20
- Signed +20, -20

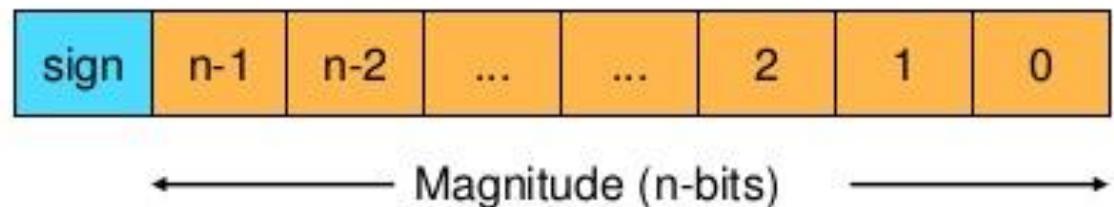
- Non-integers

- Floating point numbers - 10.25, 3.33333..., $1/8 = 0.125$

Signed Integers

- We need a way to represent negative values
- 3 representations
 - Sign & Magnitude representation (S&M)
 - Complement method
 - Bias notation or Excess notation

1. Sign & Magnitude Representation



- n -bit unsigned magnitude & sign bit (S)
- If S
 - 0 – Integer is positive or zero
 - 1 – Integer is negative or zero
- Range – $(2^n - 1)$ to $+ (2^n - 1)$

2. Complement Method

□ Base = Radix

- Radix r system means r number of symbols
- e.g., binary numbers have symbols 0, 1

□ 2 types

- r 's complement
- $(r - 1)$'s complement
- Where r is radix (base) of number system

□ Examples

- | | |
|-----------|-----------------------|
| ■ Decimal | 9's & 10's complement |
| ■ Binary | 1's & 2's complement |

Binary addition:-

A	B	Sum	Carry
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1

Binary subtraction:-

A	B	Difference	Borrow
0	0	0	0
0	1	1	1
1	0	1	0
1	1	0	0

1's complement

- Calculated by

- $(2^n - 1) - m$

- If $m = 0101$

- 1's complement of m on a 4-bit system

$$\begin{array}{r} 1 \ 1 \ 1 \ 1 \\ 0 \ 1 \ 0 \ 1 - \\ \hline 1 \ 0 \ 1 \ 0 \end{array}$$

- This represents -5 in 1's complement

Finding 1's Complement – Short Cut

- Invert each bit of m

- Example

- $m =$

- 0 0 1 0 1 0 1 1

- 1's complement of $m =$

- 1 1 0 1 0 1 0 0

- $m =$

- 0 0 0 0 0 0 0 0 = 0

- 1's complement of $m =$

- 1 1 1 1 1 1 1 1 = -0

- Values range from -127 to +127

Addition with 1's Complement

- If results has a carry add it to LSB (Least Significant Bit)
- Example

- Add 6 and -3 on a 3-bit system

- $$\begin{array}{rcl} 6 & = & 110 \\ -3 & = & \underline{100} + \\ & = & 1010 \\ & & \swarrow 1 + \\ & & 011 \end{array}$$

2's Complement

- Doesn't require end-around carry operation as in 1's complement
- 2's complement is formed by
 - Finding 1's complement
 - Add 1 to LSB
- New range is from -128 to +127
 - -128 because of +1 to negative value

Example – 2's Complement

- Find 2's complement of 0101011

- $m = 0101011$

- $= 1010100$

- $1 +$

- $2's = 1010101$

- Short-cut

1. Search for the 1st bit with value 1 starting from LSB
2. Inver all bits after 1st one

Example – 2's Complement (Cont.)

- Add 6 and -5 on a 4-bit system

$$5 = 0101$$

$$-5 = 1011$$

$$6 = 0110$$

$$\begin{array}{r} -5 \\ = \end{array} \quad \begin{array}{r} 1011 \\ + \\ 0001 \\ \hline 0001 \end{array}$$

Discard

1's vs. 2's Complement

- 1's complement has 2 zeros (+0, -0)
- Value range is less than 2's complement

- 2's complement only has a single zero
- Value range is unequal

- No need of a separate subtract circuit
 - Doing a NOT operation is much more cost effective in terms of circuit design
- However, multiplication & division is slow

Booths Multiplication Algorithm

- 2'S complement of number
- Addition or subtraction using 2's complement
- Arithmetic right shift operation
- Booths multiplication algorithm flowchart
- Booths multiplication algorithm Examples

0 1 1 0 1 1 1 0

Original binary value

1 0 0 1 0 0 0 1

1's complement

1 0 0 1 0 0 0 1	
+	1

1 0 0 1 0 0 1 0

2's complement

Binary representation of 5 is: 0101

1's Complement of 5 is: 1010

2's Complement of 5 is: (1's Complement + 1) i.e.

1010 (1's Compliment)

+ 1

1011 (2's Complement i.e. -5)

Binary representation of 3 is: 0 0 1 1

1's Complement of 3 is: 1 1 0 0

2's Complement of 3 is: (1's Complement + 1) i.e.

1 1 0 0 (1's Compliment)

$$\begin{array}{r} + 1 \\ \hline \underline{1 \ 1 \ 0 \ 1} \end{array} \text{ (2's Complement i.e. -3)}$$

Now $2 + (-3) = 0 0 1 0$ (2 in binary)

+ 1 1 0 1 (-3)

$$\underline{\underline{1 \ 1 \ 1 \ 1}} \text{ (-1)}$$

Now, to check whether it is -1 or not simply. takes 2's

Complement of -1 and kept -ve sign as it is.

-1 = 1 1 1 1

2's Complement = -(0 0 0 0)

$$\begin{array}{r} + 1 \\ \hline \underline{\underline{-(0 \ 0 \ 0 \ 1)}} \text{ i.e. -1} \end{array}$$

$$\begin{array}{r} 5 = 0\ 0\ 0\ 0\ 0\ 1\ 0\ 1 \\ \downarrow \quad \downarrow \\ 1\ 1\ 1\ 1\ 1\ 0\ 1\ 0 \\ + 1 \\ \hline - 5 = 1\ 1\ 1\ 1\ 1\ 0\ 1\ 1 \end{array}$$

Complement Digits

Add 1

$$\begin{array}{r} -13 = 1\ 1\ 1\ 1\ 0\ 0\ 1\ 1 \\ \downarrow \quad \downarrow \\ 0\ 0\ 0\ 0\ 1\ 1\ 0\ 0 \\ + 1 \\ \hline 13 = 0\ 0\ 0\ 1\ 1\ 0\ 1 \end{array}$$

Complement Digits

Add 1

2s Complement Addition/Subtraction

- Examples: 4-bit binary system

$$\begin{array}{r} +3 \\ + 4 \\ \hline \end{array} \quad \begin{array}{r} 0011 \\ + 0100 \\ \hline \end{array}$$
$$\begin{array}{r} +7 \\ \hline \end{array} \quad \begin{array}{r} 0111 \\ \hline \end{array}$$

$$\begin{array}{r} -2 \\ + -6 \\ \hline \end{array} \quad \begin{array}{r} 1110 \\ + 1010 \\ \hline \end{array}$$
$$\begin{array}{r} -8 \\ \hline \end{array} \quad \begin{array}{r} 11000 \\ \hline \end{array}$$

$$\begin{array}{r} +6 \\ + -3 \\ \hline \end{array} \quad \begin{array}{r} 0110 \\ + 1101 \\ \hline \end{array}$$
$$\begin{array}{r} +3 \\ \hline \end{array} \quad \begin{array}{r} 10011 \\ \hline \end{array}$$

$$\begin{array}{r} +4 \\ + -7 \\ \hline \end{array} \quad \begin{array}{r} 0100 \\ + 1001 \\ \hline \end{array}$$
$$\begin{array}{r} -3 \\ \hline \end{array} \quad \begin{array}{r} 1101 \\ \hline \end{array}$$



- Which of the above is/are overflow(s)?

2s complement addition and subtraction

- Simple addition and subtraction
 - simple scheme makes 2s complement the virtually unanimous choice for integer number systems in computers

$$\begin{array}{r} 4 \quad 0100 \\ + 3 \quad \underline{0011} \\ \hline 7 \quad 0111 \end{array} \qquad \begin{array}{r} -4 \quad 1100 \\ + (-3) \quad \underline{1101} \\ \hline -7 \quad 11001 \end{array}$$

$$\begin{array}{r} 4 \quad 0100 \\ - 3 \quad \underline{1101} \\ \hline 1 \quad 10001 \end{array} \qquad \begin{array}{r} -4 \quad 1100 \\ + 3 \quad \underline{0011} \\ \hline -1 \quad 1111 \end{array}$$

$$\begin{array}{r}
 1001 = -7 \\
 +\underline{0101} = 5 \\
 1110 = -2
 \end{array}$$

(a) $(-7) + (+5)$

$$\begin{array}{r}
 1100 = -4 \\
 +\underline{0100} = 4 \\
 \underline{10000} = 0
 \end{array}$$

(b) $(-4) + (+4)$

$$\begin{array}{r}
 0011 = 3 \\
 +\underline{0100} = 4 \\
 0111 = 7
 \end{array}$$

(c) $(+3) + (+4)$

$$\begin{array}{r}
 1100 = -4 \\
 +\underline{1111} = -1 \\
 \underline{11011} = -5
 \end{array}$$

(d) $(-4) + (-1)$

$$\begin{array}{r}
 0101 = 5 \\
 +\underline{0100} = 4 \\
 1001 = \text{Overflow}
 \end{array}$$

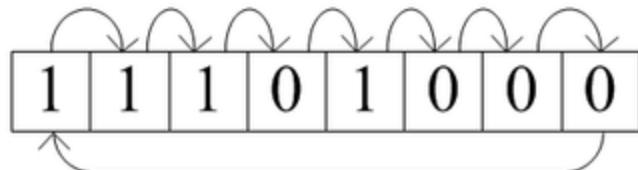
(e) $(+5) + (+4)$

$$\begin{array}{r}
 1001 = -7 \\
 +\underline{1010} = -6 \\
 \underline{10011} = \text{Overflow}
 \end{array}$$

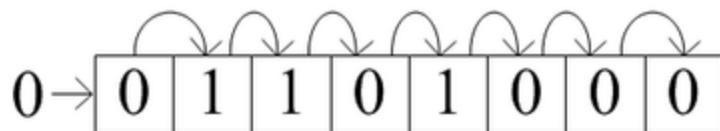
(f) $(-7) + (-6)$

Figure 9.3 Addition of Numbers in Twos Complement Representation

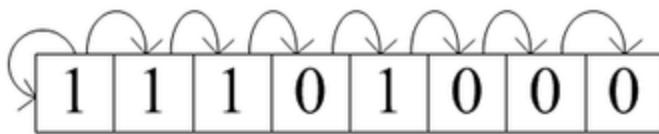
1	1	0	1	0	0	0	1
D_7	D_6	D_5	D_4	D_3	D_2	D_1	D_0



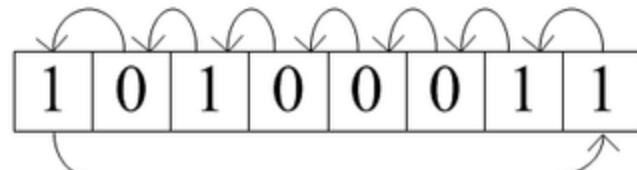
Rotate Right



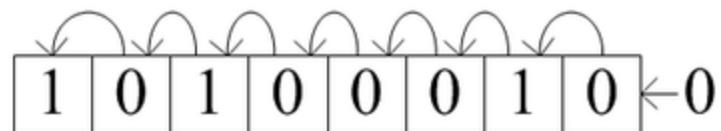
Shift Logical Right



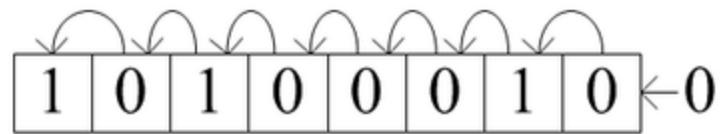
Shift Arithmetic Right



Rotate Left

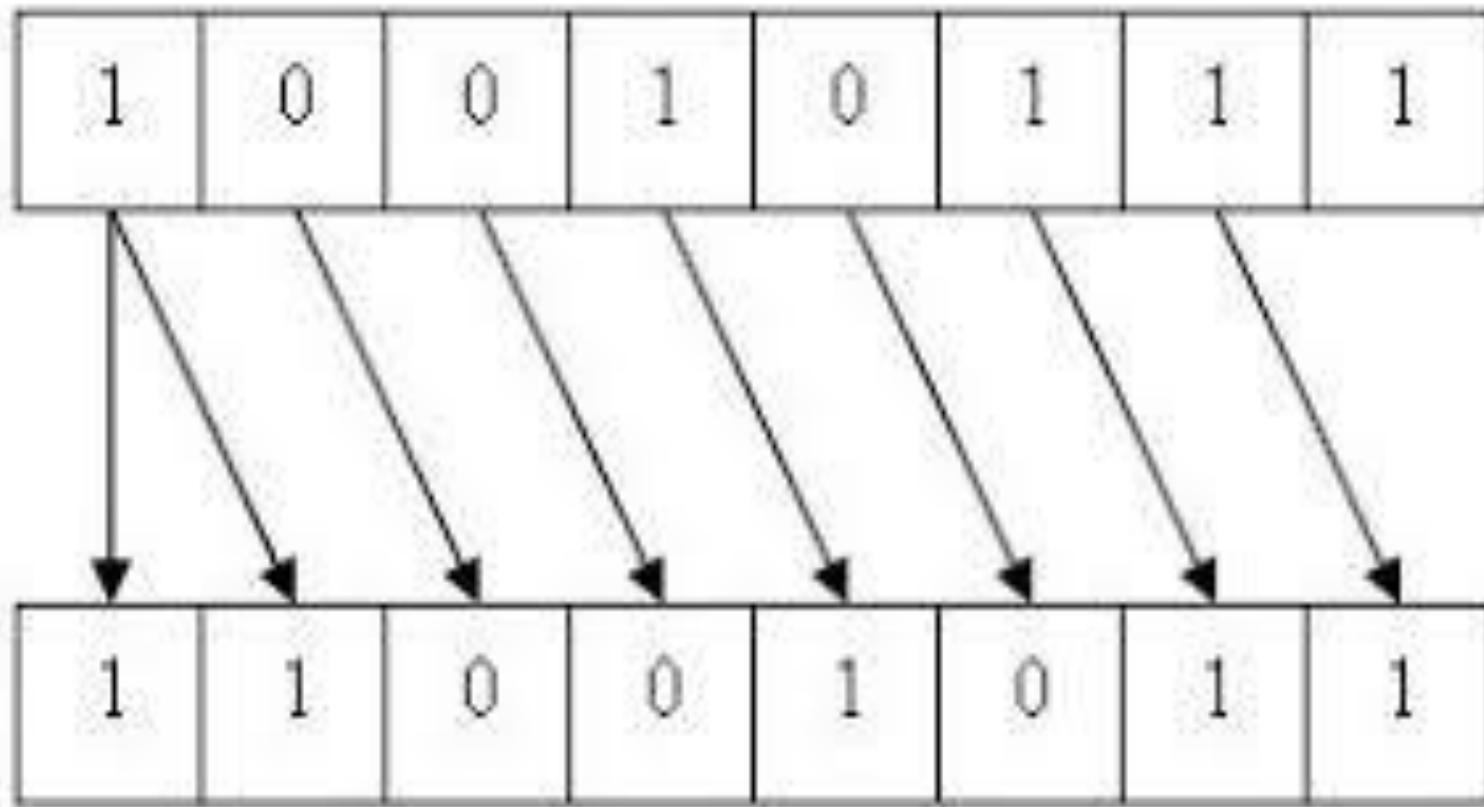


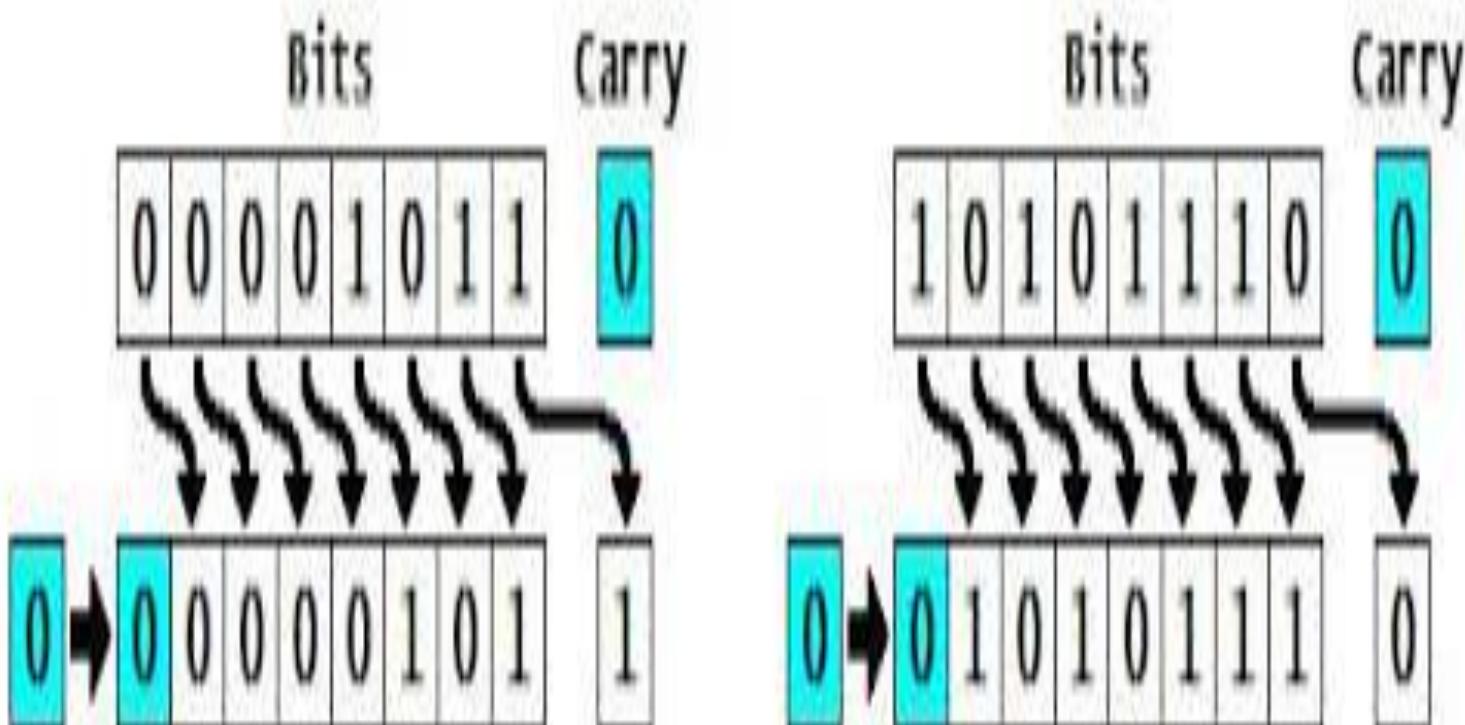
Shift Logical Left



Shift Arithmetic Left

Arithmetic's right shift

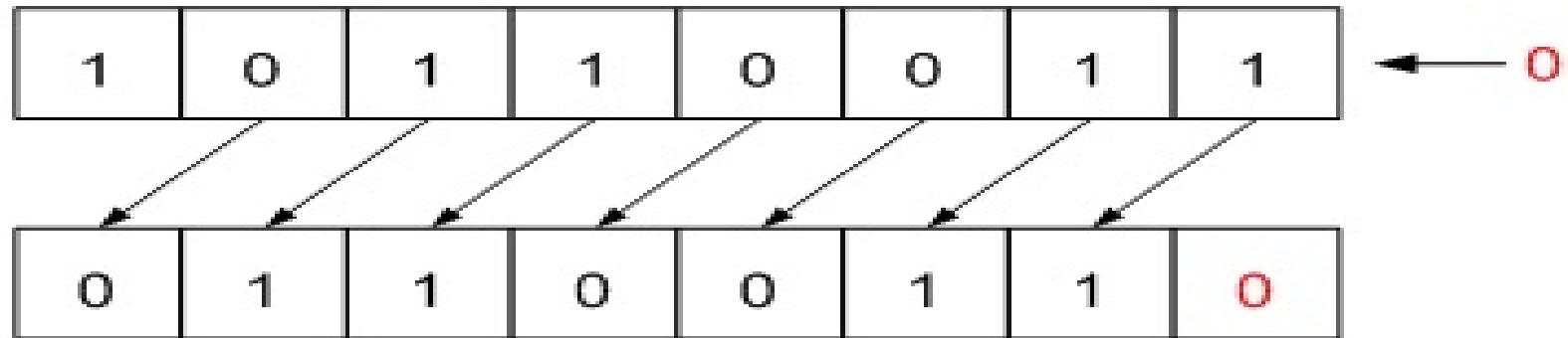




Left Logical Shift

MSB

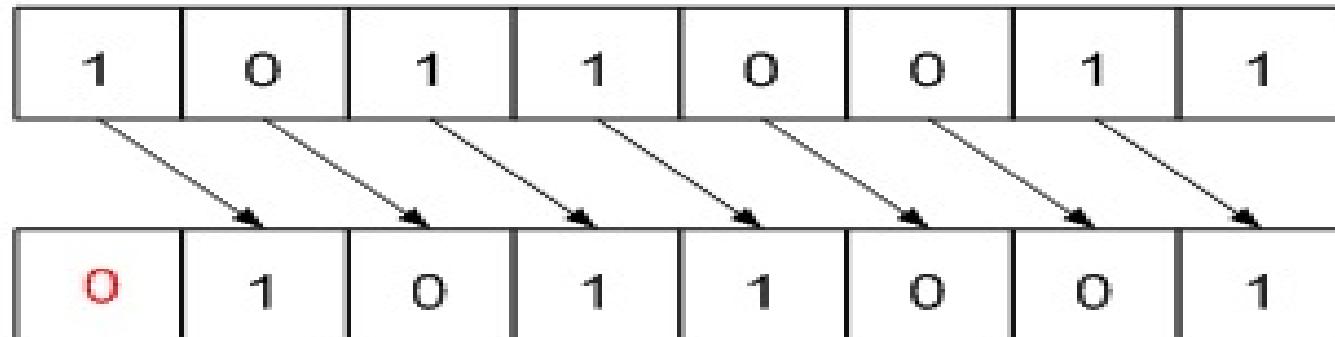
LSB



Right Logical Shift

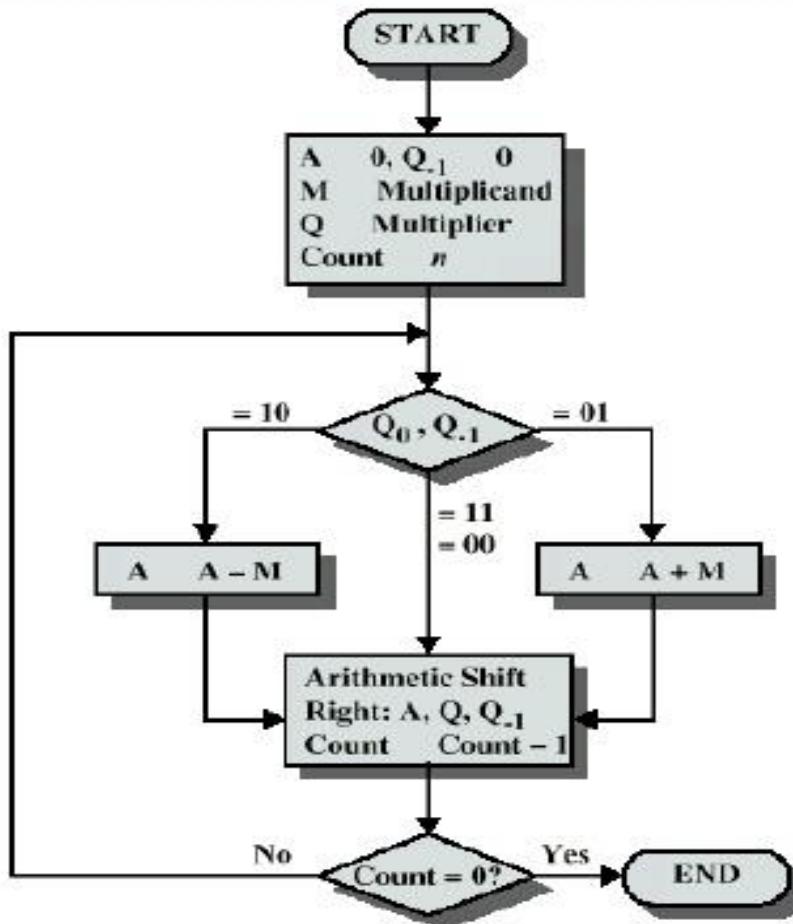
MSB

LSB

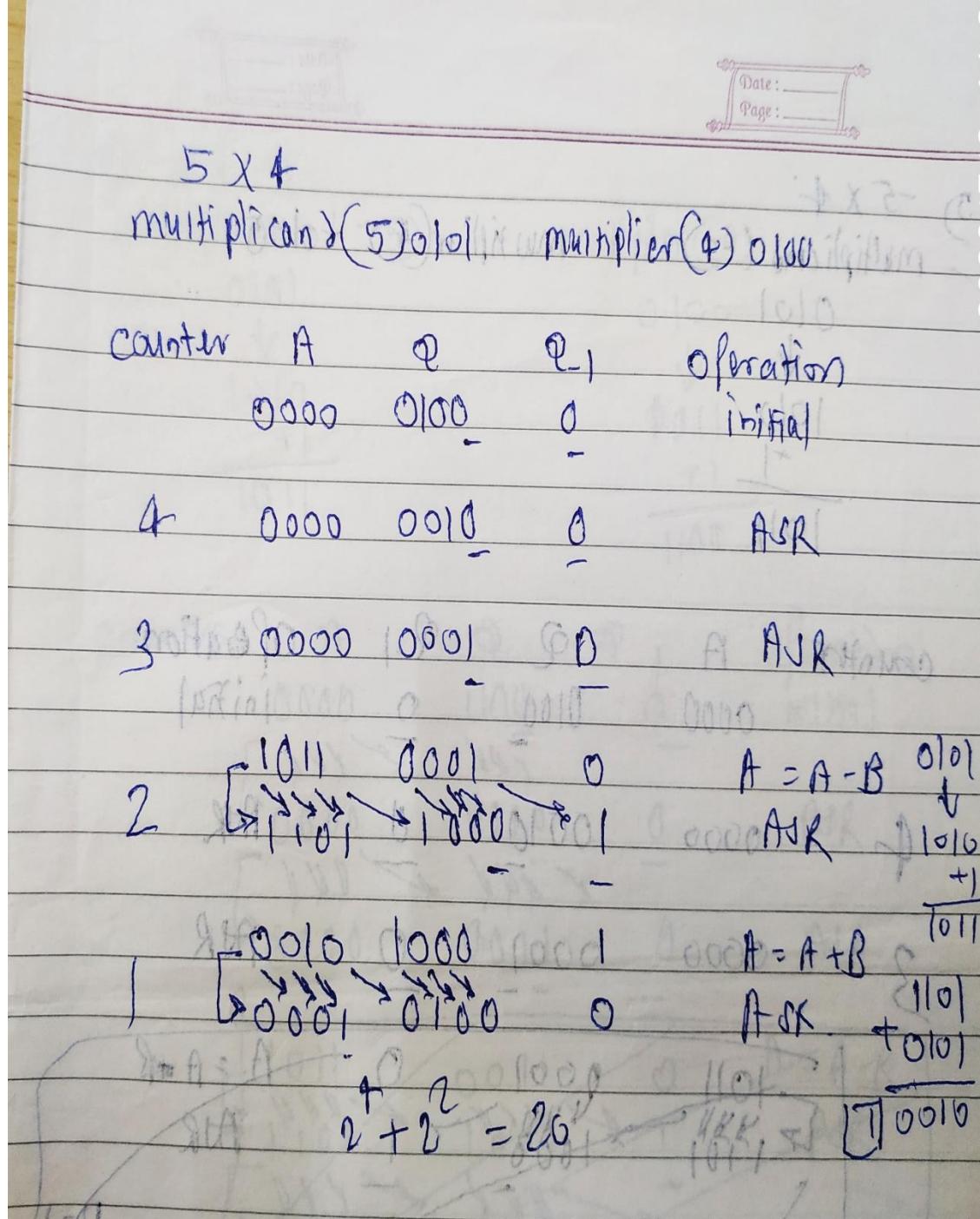
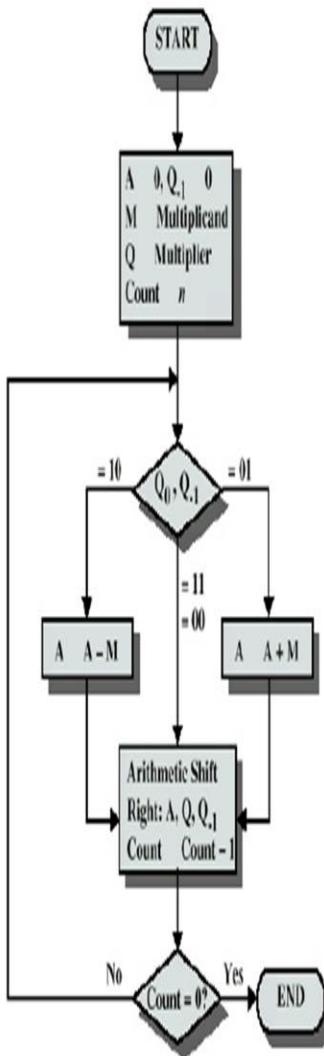


Booth's Algorithm flowchart

Booth's Algorithm



Booth's Algorithm



2s Complement Addition/Subtraction

- Examples: 4-bit binary system

$$\begin{array}{r} +3 \quad 0011 \\ + +4 \quad + 0100 \\ \hline - \quad - \quad \hline +7 \quad 0111 \\ \hline - \quad - \end{array}$$

$$\begin{array}{r} -2 \quad 1110 \\ + -6 \quad + 1010 \\ \hline - \quad - \quad \hline -8 \quad 11000 \\ \hline - \quad - \end{array}$$

$$\begin{array}{r} +6 \quad 0110 \\ + -3 \quad + 1101 \\ \hline - \quad - \quad \hline +3 \quad 10011 \\ \hline - \quad - \end{array}$$

$$\begin{array}{r} +4 \quad 0100 \\ + -7 \quad + 1001 \\ \hline - \quad - \quad \hline -3 \quad 1101 \\ \hline - \quad - \end{array}$$



- Which of the above is/are overflow(s)?

2s complement addition and subtraction

- Simple addition and subtraction
 - simple scheme makes 2s complement the virtually unanimous choice for integer number systems in computers

$$\begin{array}{r} 4 \quad 0100 \\ + 3 \quad \underline{0011} \\ \hline 7 \quad 0111 \end{array} \qquad \begin{array}{r} -4 \quad 1100 \\ + (-3) \quad \underline{1101} \\ \hline -7 \quad 11001 \end{array}$$

$$\begin{array}{r} 4 \quad 0100 \\ - 3 \quad \underline{1101} \\ \hline 1 \quad 10001 \end{array} \qquad \begin{array}{r} -4 \quad 1100 \\ + 3 \quad \underline{0011} \\ \hline -1 \quad 1111 \end{array}$$

$$\begin{array}{r}
 1001 = -7 \\
 +\underline{0101} = 5 \\
 1110 = -2
 \end{array}$$

(a) $(-7) + (+5)$

$$\begin{array}{r}
 1100 = -4 \\
 +\underline{0100} = 4 \\
 \underline{10000} = 0
 \end{array}$$

(b) $(-4) + (+4)$

$$\begin{array}{r}
 0011 = 3 \\
 +\underline{0100} = 4 \\
 0111 = 7
 \end{array}$$

(c) $(+3) + (+4)$

$$\begin{array}{r}
 1100 = -4 \\
 +\underline{1111} = -1 \\
 \underline{11011} = -5
 \end{array}$$

(d) $(-4) + (-1)$

$$\begin{array}{r}
 0101 = 5 \\
 +\underline{0100} = 4 \\
 1001 = \text{Overflow}
 \end{array}$$

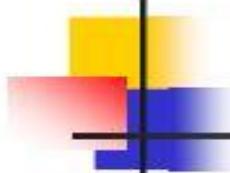
(e) $(+5) + (+4)$

$$\begin{array}{r}
 1001 = -7 \\
 +\underline{1010} = -6 \\
 \underline{10011} = \text{Overflow}
 \end{array}$$

(f) $(-7) + (-6)$

Figure 9.3 Addition of Numbers in Twos Complement Representation

Booths multiplication

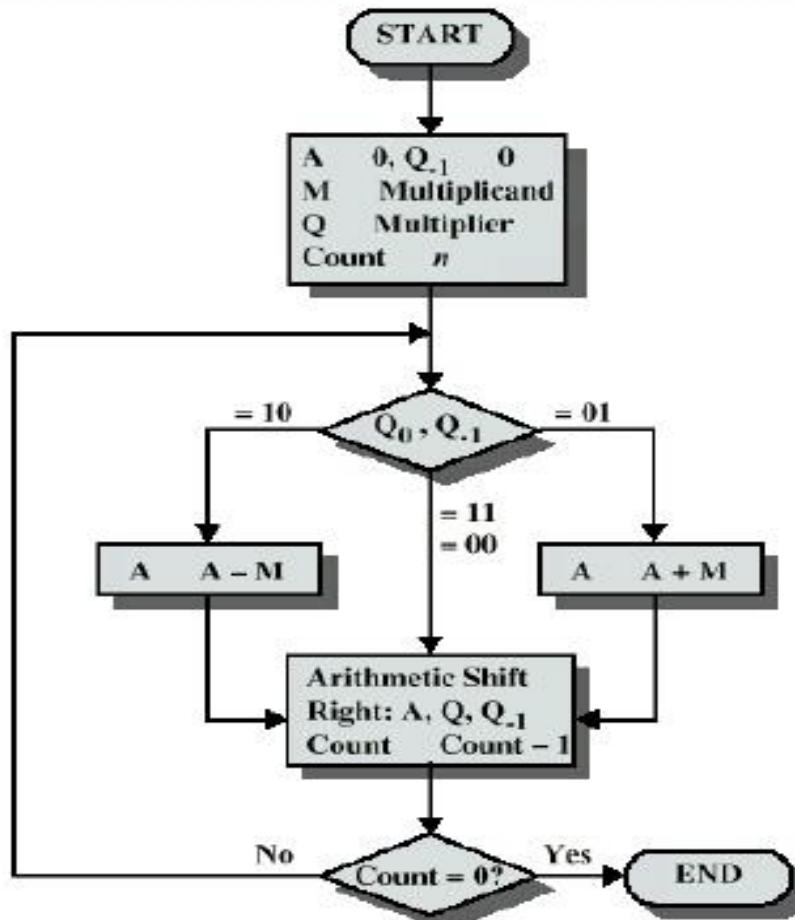


Points to remember

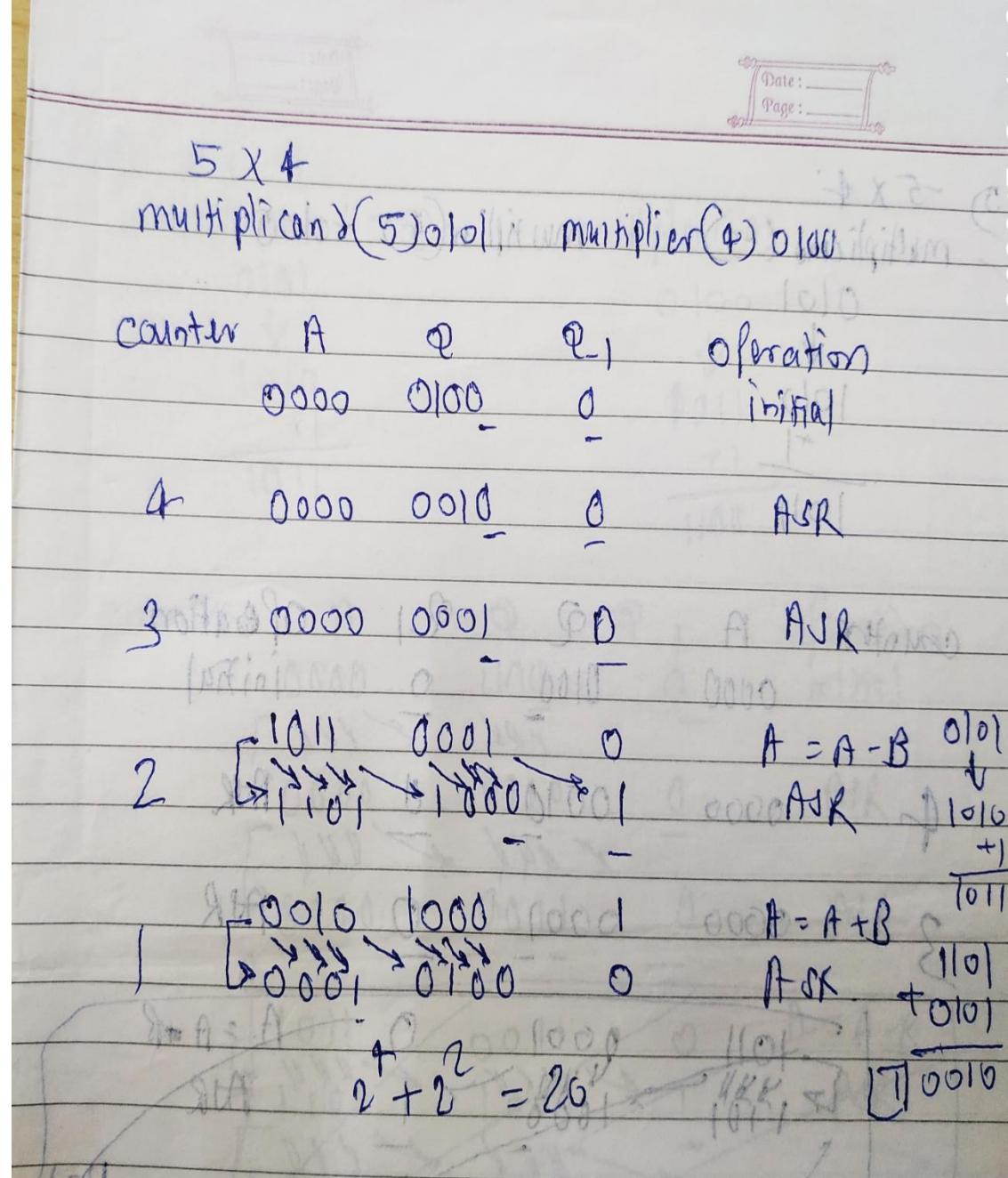
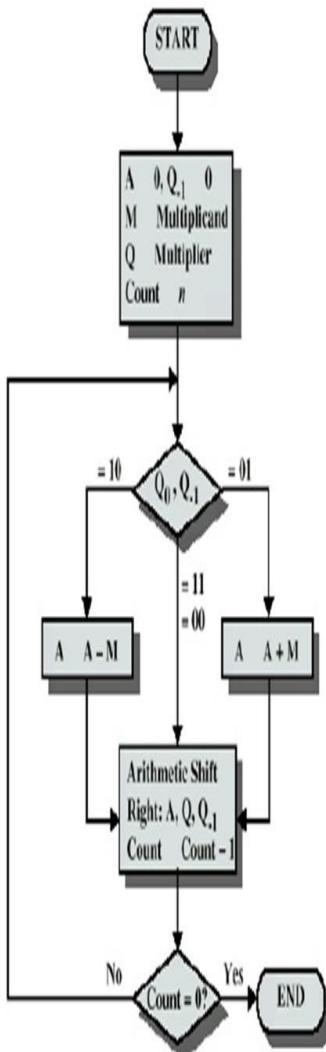
- When using Booth's Algorithm:
 - You will need twice as many bits in your **product** as you have in your original two **operands**.
 - The **leftmost bit** of your operands (both your multiplicand and multiplier) is a **SIGN** bit, and cannot be used as part of the value.

Booth's Algorithm flowchart

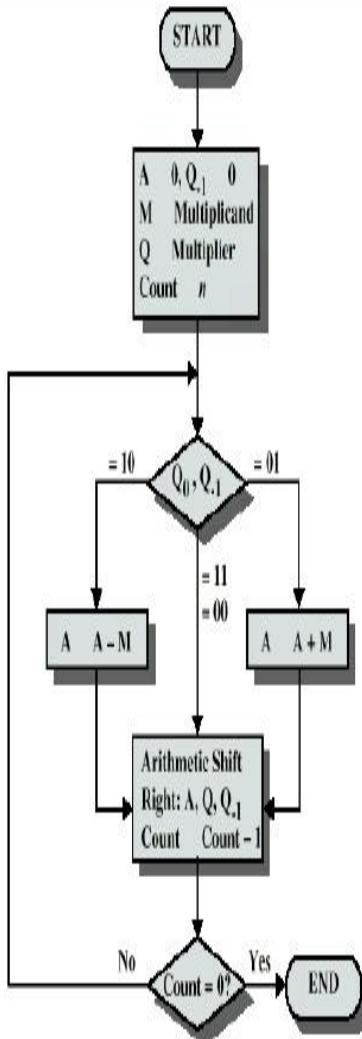
Booth's Algorithm



Booth's Algorithm



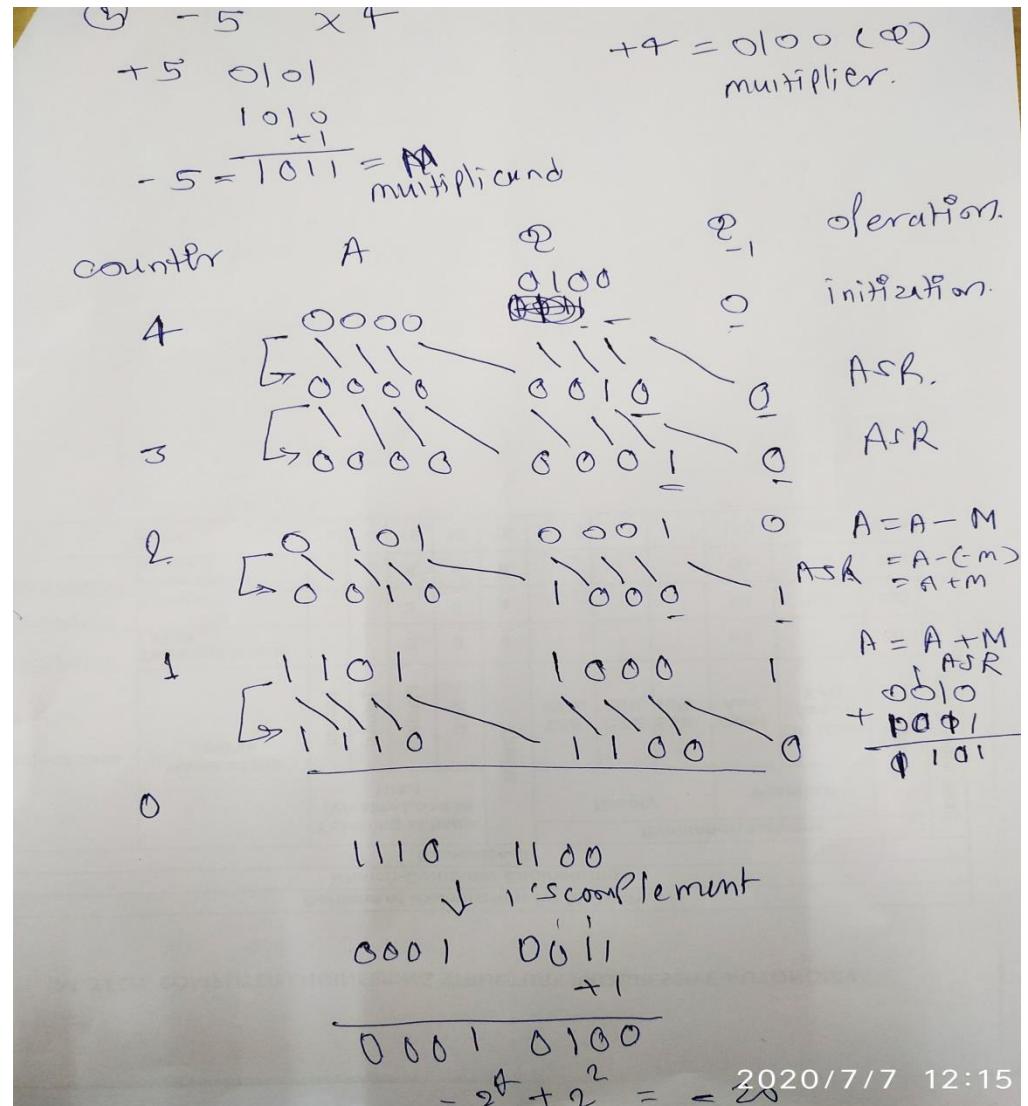
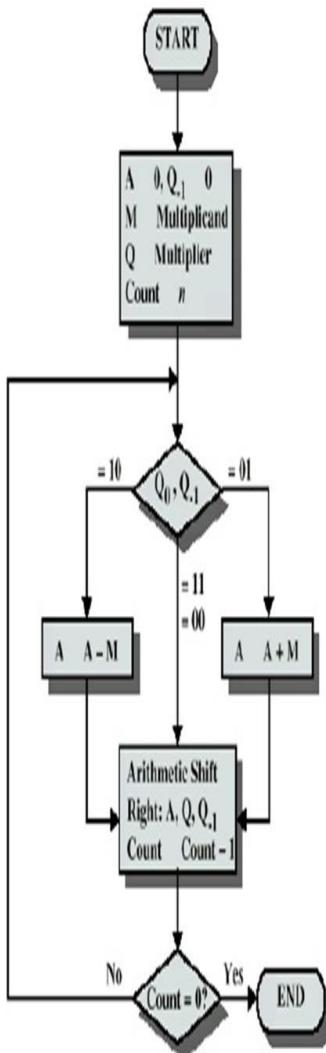
Booth's Algorithm



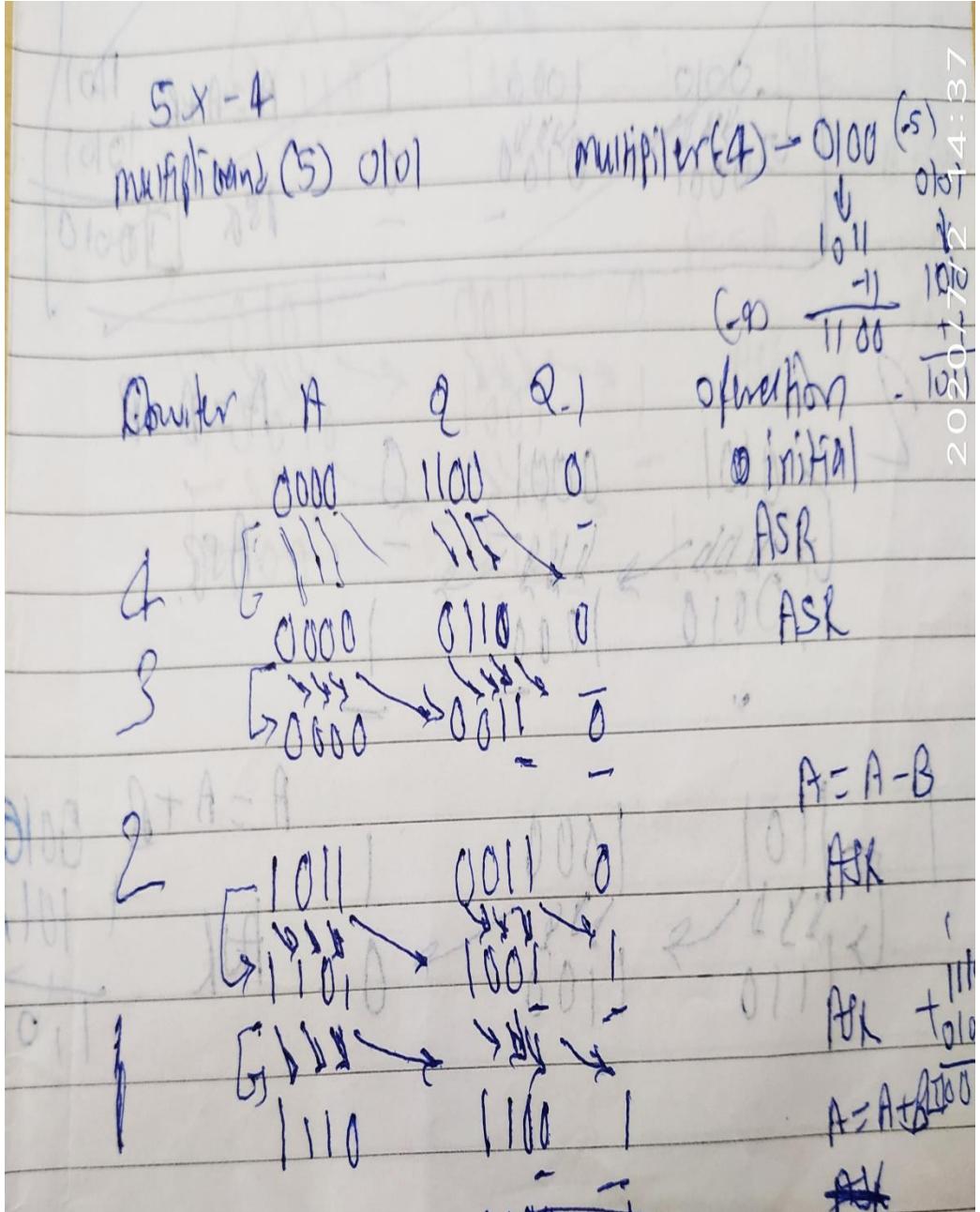
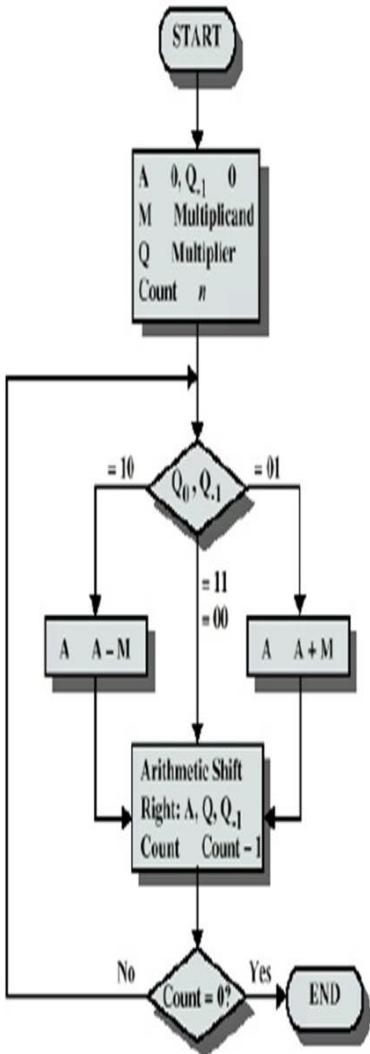
Example of Booth's Algorithm

A	Q	Q_{-1}	M	
0000	0011	0	0111	Initial Values
1001	0011	0	0111	A A - M } First
1100	1001	1	0111	Shift } Cycle
1110	0100	1	0111	Shift } Second
0101	0100	1	0111	A A + M } Cycle
0010	1010	0	0111	Shift } Third
0001	0101	0	0111	Shift } Fourth

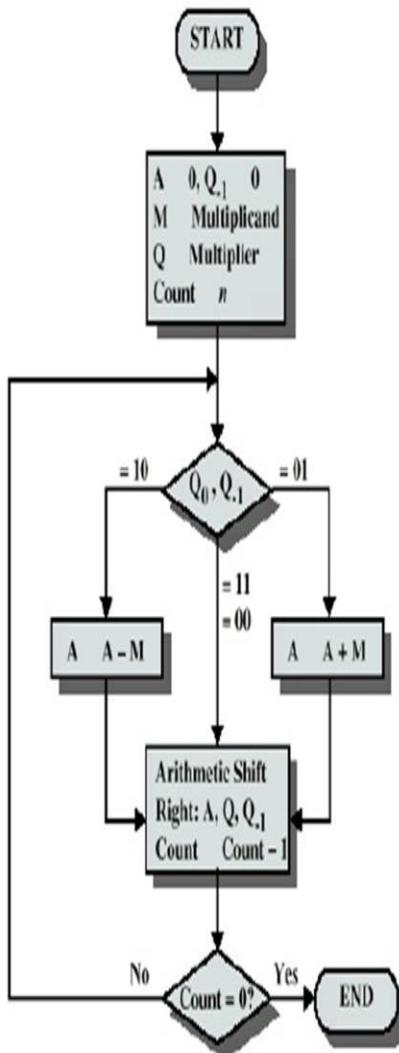
Booth's Algorithm



Booth's Algorithm

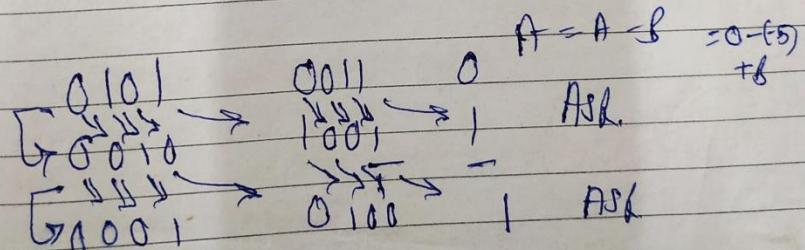
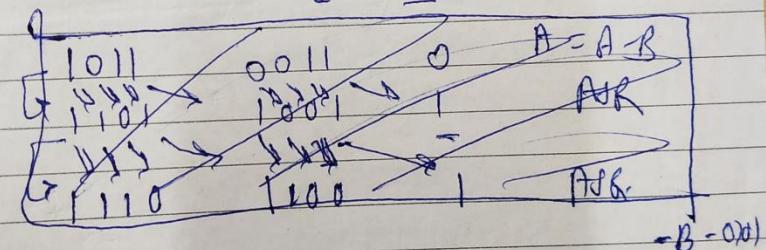


Booth's Algorithm

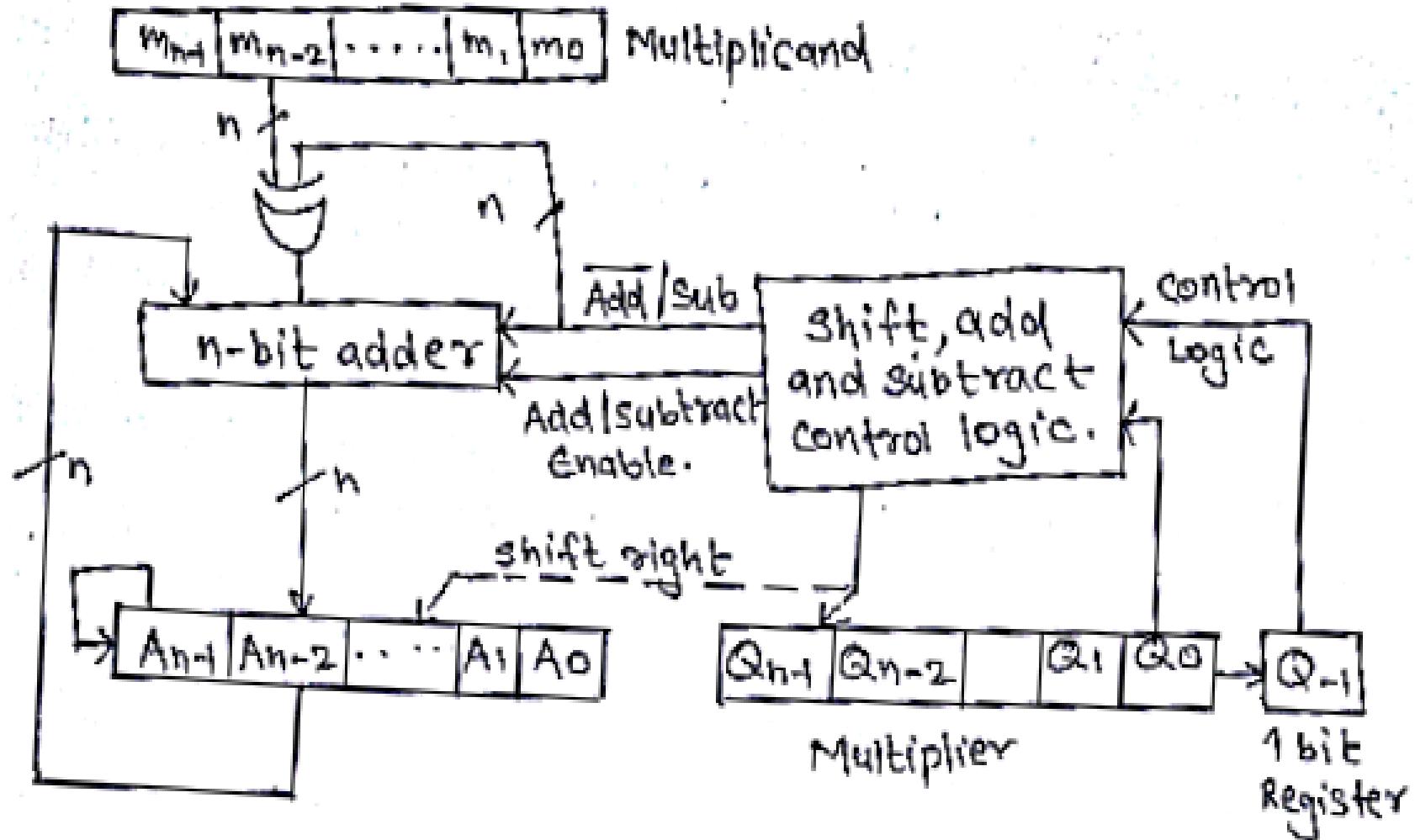


$$\begin{array}{r}
 \text{1) } -5 \times -4 \\
 \text{multiplicand } (-5) \qquad \qquad \qquad \text{multiplier } (-4) \\
 0101 \\
 \downarrow \\
 1010 \\
 +1 \\
 \hline
 1011
 \end{array}$$

Counter	A	Q	Q ₋₁	operation
	0000	1100	0	
[1011]	0000	0110	0	ASR
[1011]	0000	0011	0	ASR

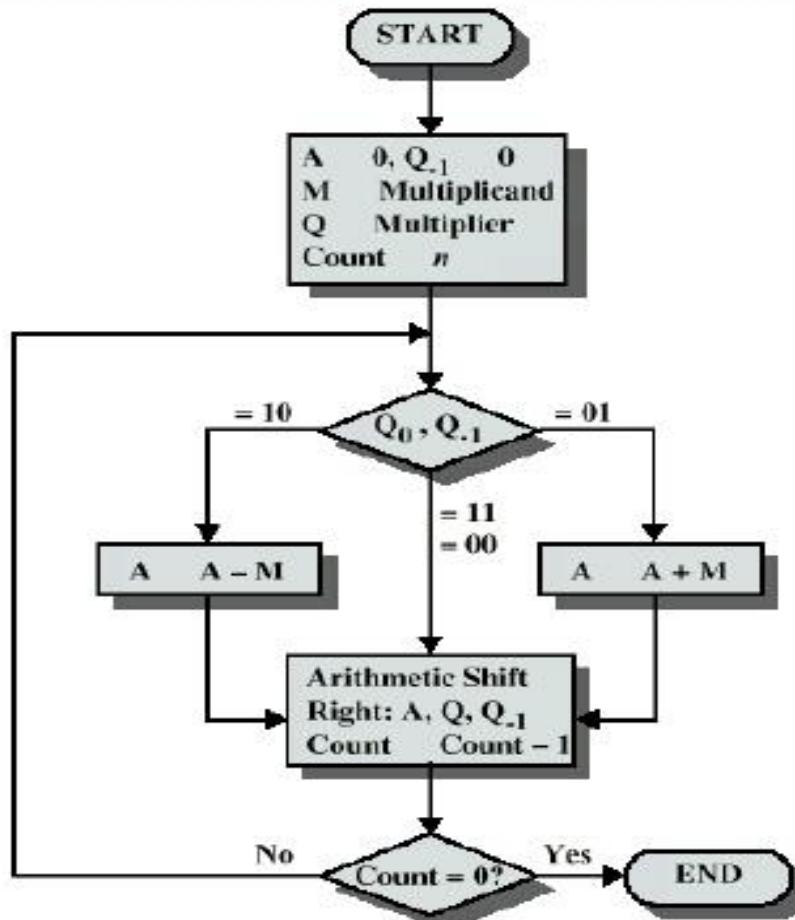


Booth's algorithm hardware implementation

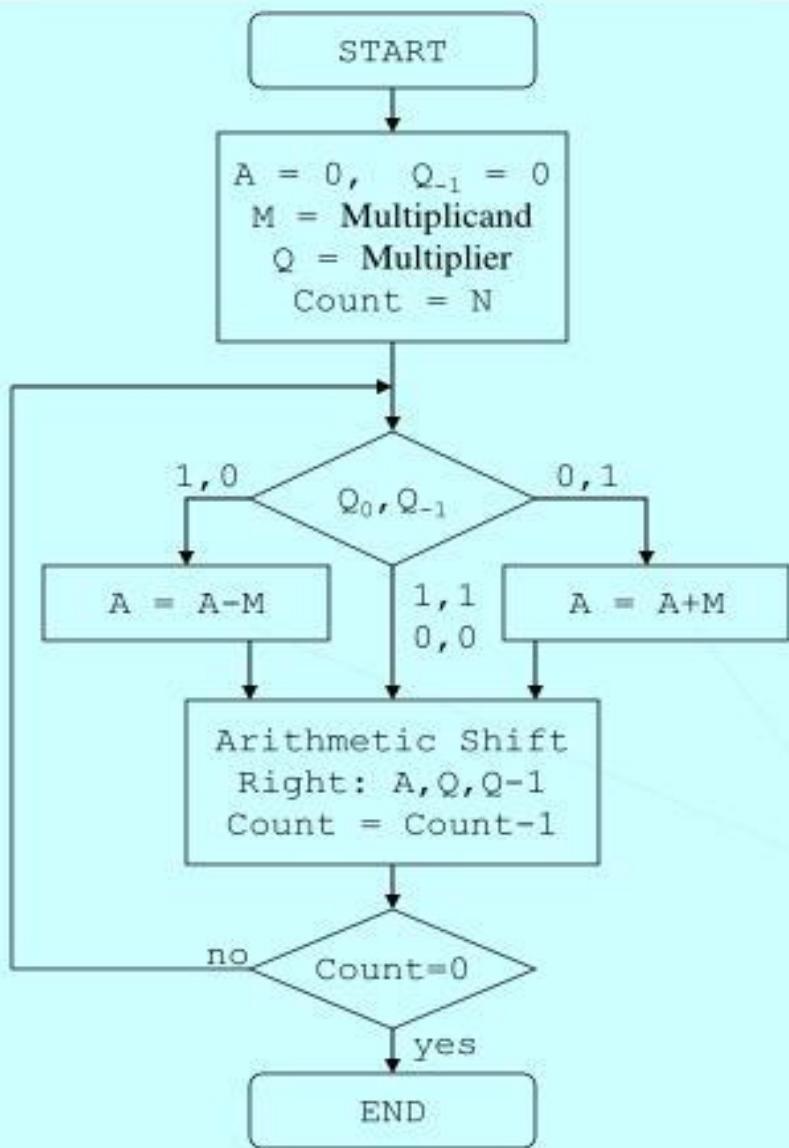


Booth's Algorithm flowchart

Booth's Algorithm



Booth's Algorithm for 2's Complement Multiplication



$$\begin{array}{l} M = 0101 \quad -M = 1011 \\ Q = 0110 \\ N = 4 \end{array}$$

A	Q	Q_{-1}	N
0000	0110	0	4
0000	0011	0	3
1011	0011	0	
1101	1001	1	2
1110	1100	1	1
0011	1100	1	
0001	1110	0	0

Answer is in A, Q

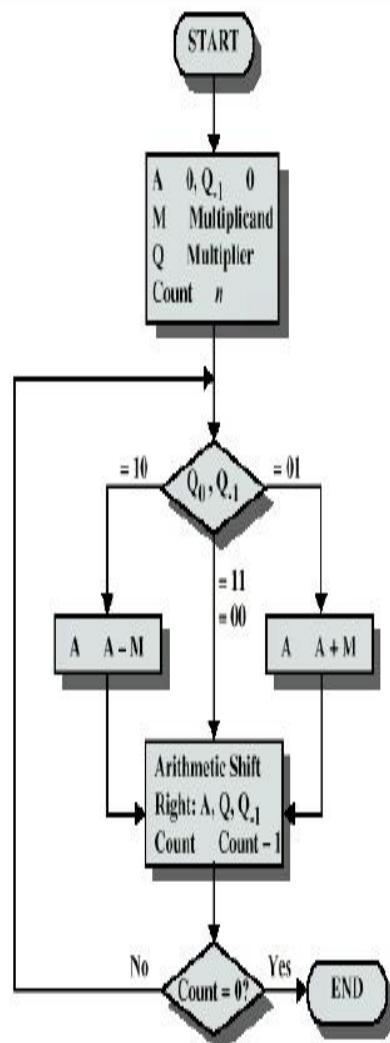
$$00011110_2 = 30_{10}$$

Two's complement multiplication can be performed using Booth's Algorithm.

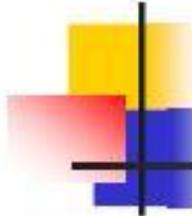
To save time, both M and $-M$ can be computed ahead and provided as one of the inputs to the adders for $A-M$ and $A+M$.

Instead of a loop, the math processor can provide separate hardware for each stage of the multiplication algorithm.

Booth's Algorithm



	A	Q	Q ₁	M=0010
6 * 2 = 12	0000	0010	0	
6 = 0110	0000	0001	0	
2 = 0010				
Q ₀ Q ₁	1010	0001	0	
0 0	1101	0000	1	
1 1				
1 0	0011	0000	1	
0 1	0001	1000	0	
	0000	1100	0	



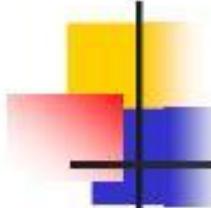
To begin

- Decide which operand will be the **multiplier** and which will be the **multiplicand**
- Convert both operands to **two's complement** representation using X bits
 - X must be at least one more bit than is required for the binary representation of the numerically larger operand
- Begin with a product that consists of the multiplier with an additional X leading zero bits



Example

- In the week by week, there is an example of multiplying **2 x (-5)**
- For our example, let's reverse the operation, and multiply **(-5) x 2**
 - The numerically larger operand (5) would require 3 bits to represent in binary (101). So we must use AT LEAST 4 bits to represent the operands, to allow for the sign bit.
- Let's use 5-bit 2's complement:
 - -5 is **11011** (multiplier)
 - 2 is **00010** (multiplicand)



Beginning Product

- The multiplier is:

11011

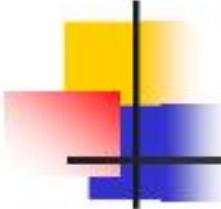
- Add 5 leading zeros to the **multiplier** to get the **beginning product**:

00000 11011



Step 1 for each pass

- Use the **LSB** (least significant bit) and the **previous LSB** to determine the arithmetic action.
 - If it is the FIRST pass, use **0** as the previous LSB.
- Possible arithmetic actions:
 - **00** → no arithmetic operation
 - **01** → add multiplicand to left half of product
 - **10** → subtract multiplicand from left half of product
 - **11** → no arithmetic operation



Step 2 for each pass

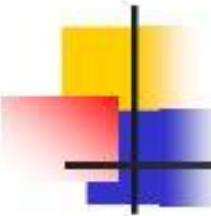
- Perform an **arithmetic right shift** (ASR) on the entire product.
- NOTE: For X-bit operands, Booth's algorithm requires X passes.



Example

- Let's continue with our example of multiplying
(-5) x 2
- Remember:
 - -5 is **11011** (multiplier)
 - 2 is **00010** (multiplicand)
- And we added 5 leading zeros to the
multiplier to get the **beginning product**:

00000 11011



Example continued

- Initial Product and previous LSB

00000 11011 0

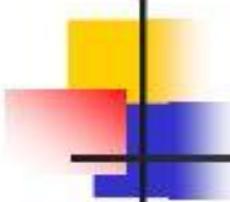
(Note: Since this is the first pass, we use 0 for the previous LSB)

- Pass 1, Step 1: Examine the last 2 bits

00000 11011 0

The last two bits are **10**, so we need to:

subtract the **multiplicand** from left half of product



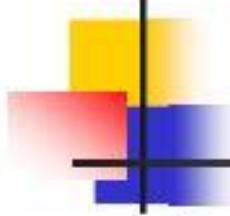
Example: Pass 1 continued

- Pass 1, Step 1: Arithmetic action

(1) $\begin{array}{r} 00000 \\ -00010 \\ \hline 11110 \end{array}$ (left half of product)
(multiplicand)
(uses a phantom borrow)

- Place result into **left half** of product

11110 11011 0



Example: Pass 1 continued

- Pass 1, Step 2: ASR (arithmetic shift right)

- Before ASR

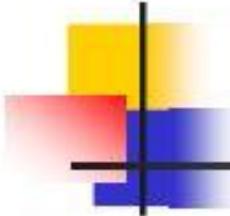
11110 11011 0

- After ASR

11111 01101 1

(left-most bit was 1, so a 1 was shifted in on the left)

- Pass 1 is complete.



Example: Pass 2

- Current Product and previous LSB

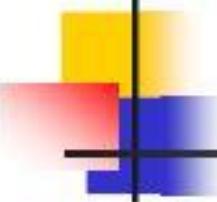
11111 01101 1

- Pass 2, Step 1: Examine the last 2 bits

11111 01101 1

The last two bits are **11**, so we do NOT need to perform an arithmetic action --

just proceed to step 2.



Example: Pass 2 continued

- Pass 2, Step 2: ASR (arithmetic shift right)

- Before ASR

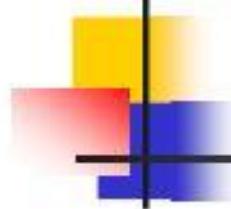
11111 01101 1

- After ASR

11111 10110 1

(left-most bit was 1, so a 1 was shifted in on the left)

- Pass 2 is complete.



Example: Pass 3

- Current Product and previous LSB

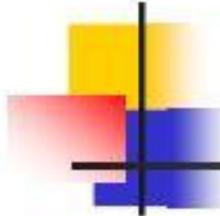
11111 10110 1

- Pass 3, Step 1: Examine the last 2 bits

11111 10110 1

The last two bits are **01**, so we need to:

add the **multiplicand** to the left half of the product



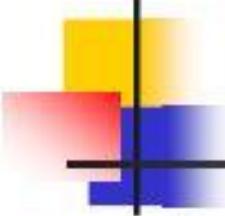
Example: Pass 3 continued

- Pass 3, Step 1: Arithmetic action

(1)
$$\begin{array}{r} 11111 \\ +00010 \\ \hline 00001 \end{array}$$
 (left half of product)
(mulitplicand)
(drop the leftmost carry)

- Place result into **left half** of product

00001 10110 1



Example: Pass 3 continued

- Pass 3, Step 2: ASR (arithmetic shift right)

- Before ASR

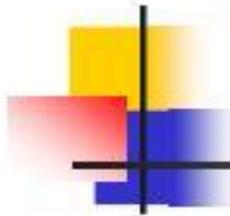
00001 10110 1

- After ASR

00000 11011 0

(left-most bit was 0, so a 0 was shifted in on the left)

- Pass 3 is complete.



Example: Pass 4

- Current Product and previous LSB

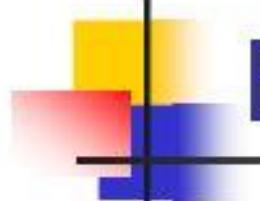
00000 11011 0

- Pass 4, Step 1: Examine the last 2 bits

00000 11011 0

The last two bits are **10**, so we need to:

subtract the **multiplicand** from the left half of the product



Example: Pass 4 continued

- Pass 4, Step 1: Arithmetic action

(1) $\begin{array}{r} 00000 \\ -00010 \\ \hline 11110 \end{array}$ (left half of product)
(multiplicand)
(uses a phantom borrow)

- Place result into **left half** of product

11110 11011 0



Example: Pass 4 continued

- Pass 4, Step 2: ASR (arithmetic shift right)

- Before ASR

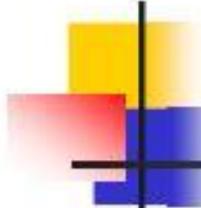
11110 11011 0

- After ASR

11111 01101 1

(left-most bit was 1, so a 1 was shifted in on the left)

- Pass 4 is complete.



Example: Pass 5

- Current Product and previous LSB

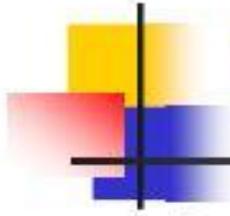
11111 01101 1

- Pass 5, Step 1: Examine the last 2 bits

11111 01101 1

The last two bits are **11**, so we do NOT need to perform an arithmetic action --

just proceed to step 2.



Example: Pass 5 continued

- Pass 5, Step 2: ASR (arithmetic shift right)

- Before ASR

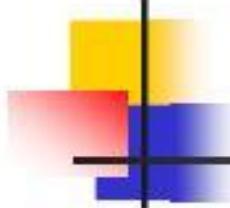
```
11111 01101 1
```

- After ASR

```
11111 10110 1
```

(left-most bit was 1, so a 1 was shifted in on the left)

- Pass 5 is complete.



Final Product

- We have completed 5 passes on the 5-bit operands, so we are done.
- Dropping the **previous** LSB, the resulting **final product** is:

11111 10110



Verification

- To confirm we have the correct answer, convert the 2's complement **final product** back to decimal.
- Final product: **11111 10110**
- Decimal value: **-10**
which is the CORRECT product of:

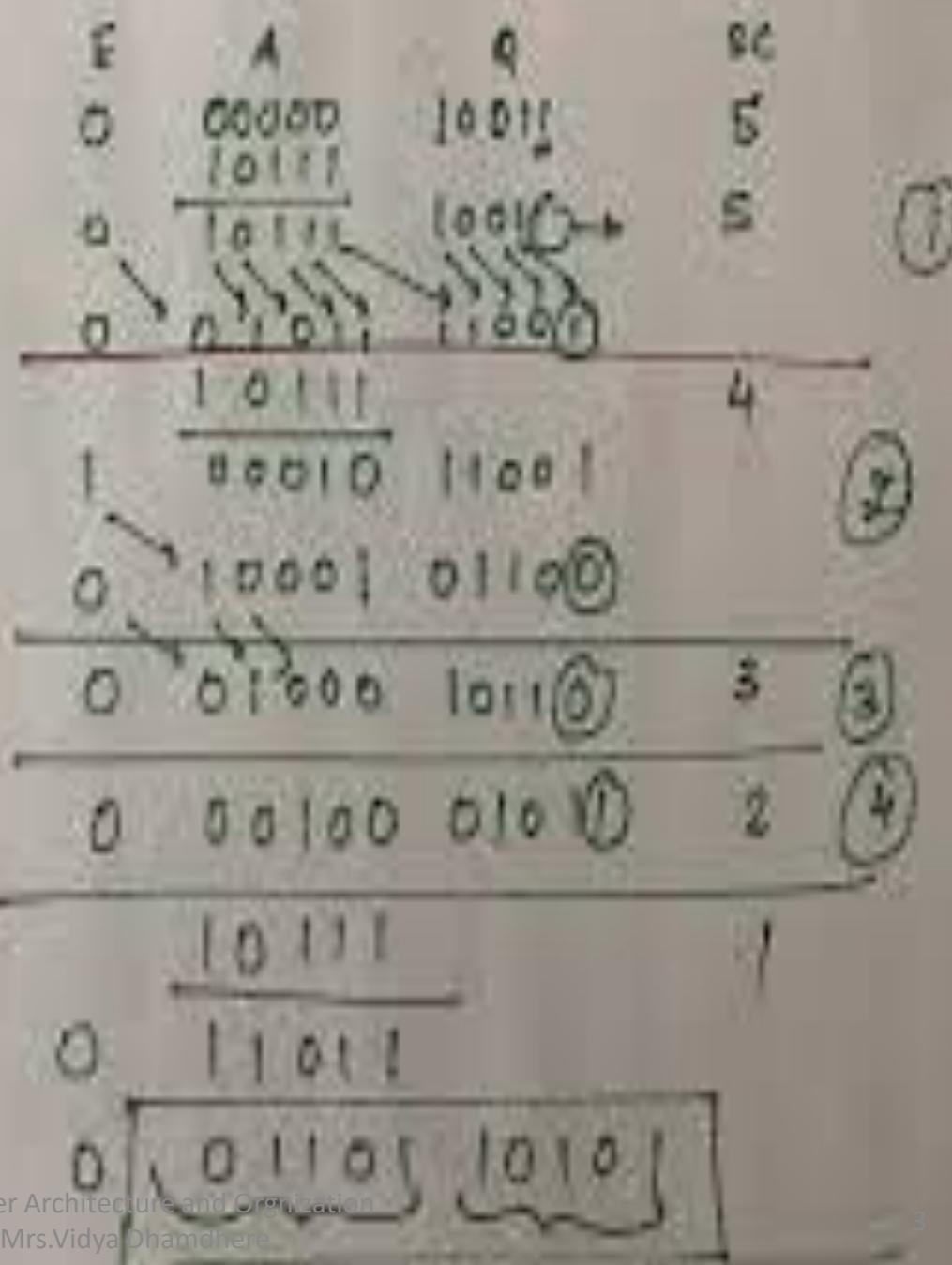
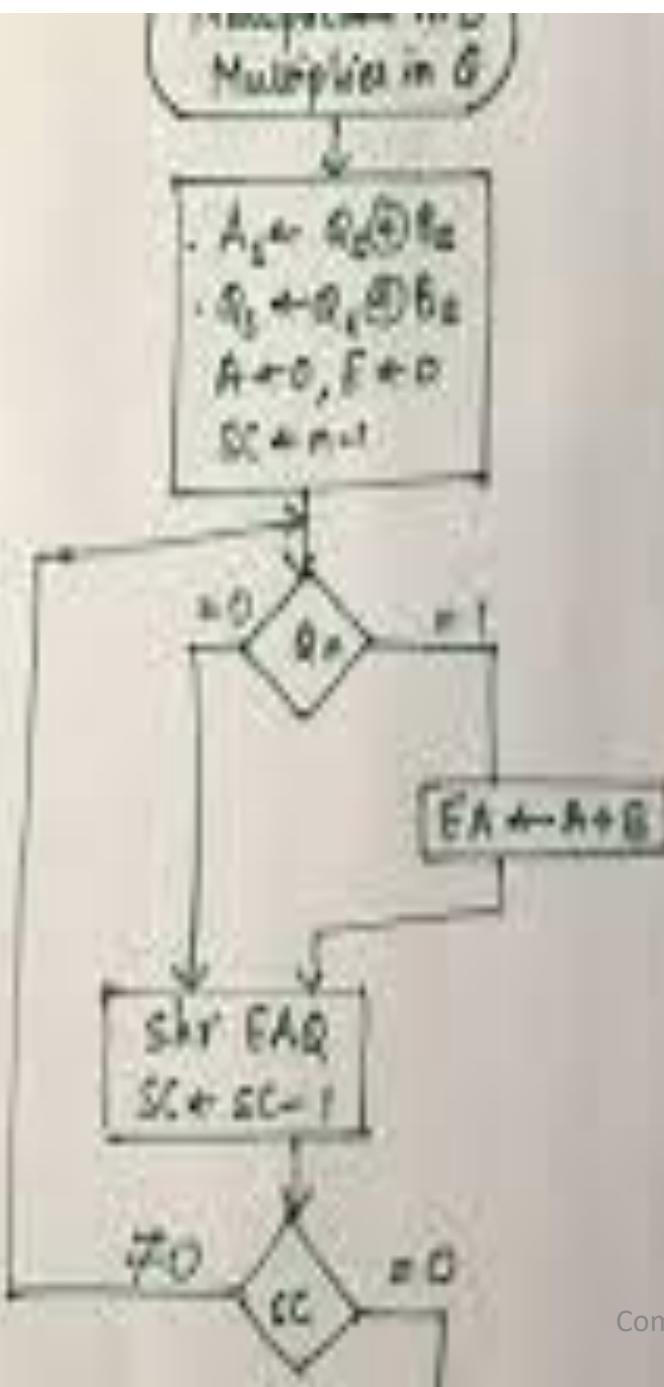
$$(-5) \times 2$$

UNIT II Data path unit

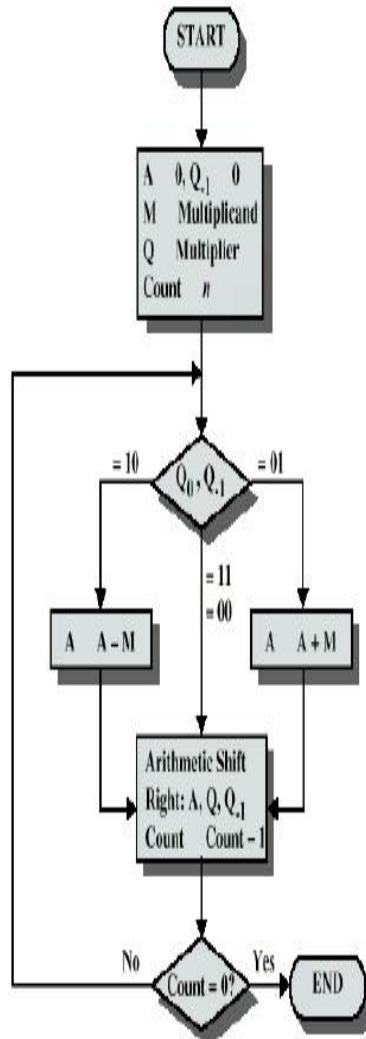
Data Representation, Fixed and Floating point numbers, Signed numbers, Fixed-Point Arithmetic, Booths Algorithm-multiplication , Arithmetic Logic unit, Floating point representations, IEEE standards, Floating point arithmetic.

Booths Multiplication Algorithm

- 2'S complement of number
- Addition or subtraction using 2's complement
- Arithmetic right shift operation
- Booths multiplication algorithm flowchart
- Booths multiplication algorithm Examples



Booth's Algorithm



Example of Booth's Algorithm

A	Q	Q-1	M	
0000	0011	0	0111	Initial Values
1001	0011	0	0111	A A - M } First
1100	1001	1	0111	Shift } Cycle
1110	0100	1	0111	Shift } Second
0101	0100	1	0111	A A + M } Third
0010	1010	0	0111	Shift } Cycle
0001	0101	0	0111	Shift } Fourth

multiplication of -5×4 ($-5 = 1011$) = B

steps

A

\otimes

\otimes^{-1}

operation

1

0000

0100

0

Initial

2.

0000

0010

0

Right shift

3.

0101

0001

0

$A \leftarrow A - B$

0010

1000

1

Right shift

4.

1101

1000

1

$A \leftarrow A + B$

1110 1100 0

Right shift

-20

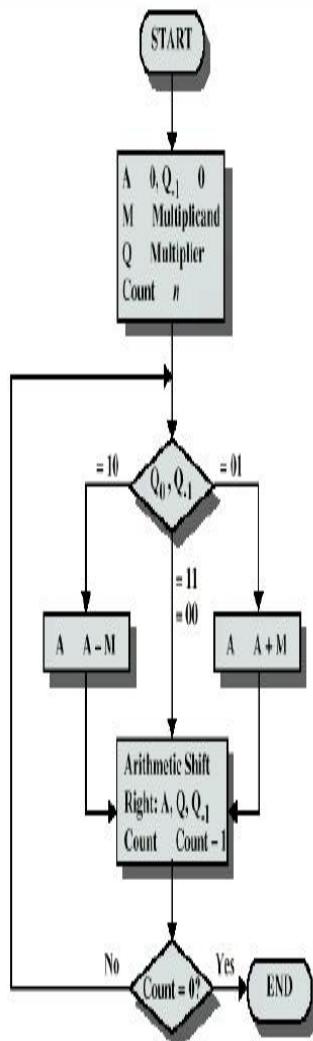
Verify:

0001 00

$-5 \times 4 \neq -20$

$5 = 0101$

Booth's Algorithm



Booth's Algorithm Example

$$\begin{array}{r}
 \text{Multiply } 4 \times -9 \\
 \begin{array}{r}
 00100 \\
 \times 10111
 \end{array}
 \end{array}$$

	LHS prod	RHS prod	
Step 0	00000	101110	(sub 4 = add -4)
Step 1	+11100		
Step 2	11100	10111	(shift after add)
Step 3	11110	01011	(shift without add)
Step 4	+00100		(add +4)
Step 5	00011	10010	
	+11100		(sub 4 = add -4)
	11101	110010	
	11110	111001	(shift after add)

Remember

$$4 = 00100$$

$$-4 = 11100$$

Feb 2005

$$1111011100 = -(0000100011 + 1)$$

$$= -(0000100100)$$

$$= -36 = 4 \times -9!$$

Multiplication and Division

19

CASE 2 – Negative Multiplier: $(5)_{10} \times (-4)_{10}$

Multiplicand (B) = 0101(5) , Multiplier (Q) = 1100(-4)

Steps	A	Q	Q-1	Operation
	0000	110 <u>0</u>	<u>0</u>	Initial
Step 1 :	0000	011 <u>0</u>	<u>0</u>	Right Shift
Step 2 :	0000	001 <u>1</u>	<u>0</u>	Right Shift
Step 3 :	1011	001 <u>1</u>	<u>0</u>	$A \leftarrow A - B$
	1101	1001	1	Right Shift
Step 4:	1110	110 <u>0</u>	<u>1</u>	Right Shift
	1110	1100		

RECORDED WITH

SCREENCAST
Result:

MATIC

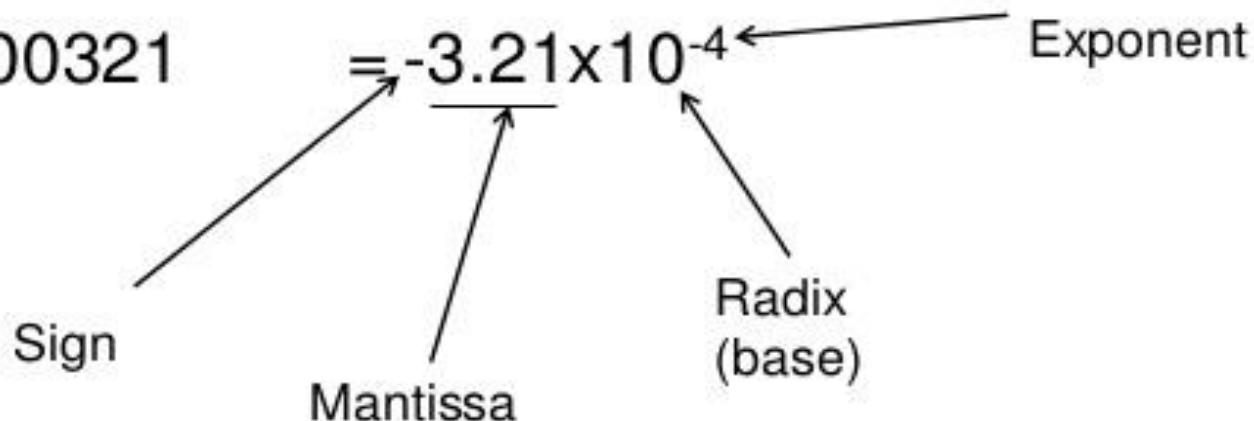
11111000 = 20 (2's complement of 20)
Mrs.Vidya Dhamdhere

Floating Point Numbers

- We needed to represent fractional values & values beyond $2^n - 1$

- $+3207.23 = 3.20723 \times 10^3$

- $-0.000321 = -3.21 \times 10^{-4}$

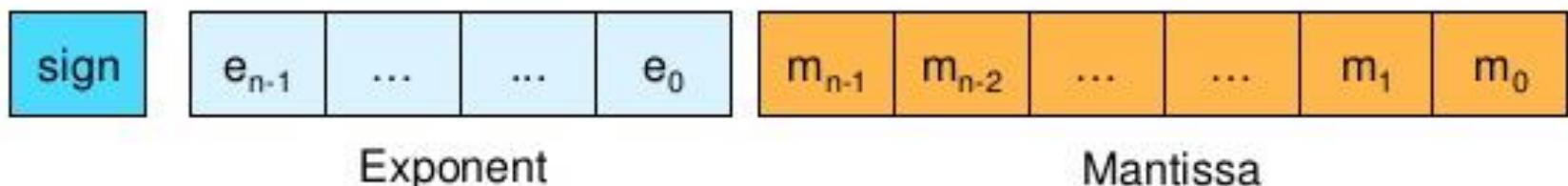


Floating Point Numbers (Cont.)

$$N = (-1)^s m \cdot 2^e$$

Diagram illustrating the components of a floating-point number N :

- Sign**: Points to the sign bit s .
- Exponent**: Points to the exponent bits e .
- Radix**: Points to the radix point (implied 1) between the sign and exponent.
- Mantissa**: Points to the fraction part m .

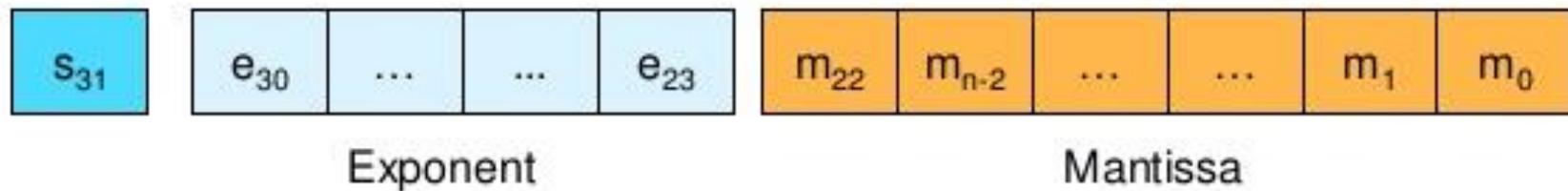


IEEE Floating Point Standard (FPS)

- 2 standards
 - 1. Single precision
 - 32-bits
 - 23-bit mantissa
 - 8-bit exponent
 - 1-bit sign
 - 2. Double precision
 - 64-bits
 - 52-bit mantissa
 - 11-bit exponent
 - 1-bit sign



IEEE Floating Point Standard (Cont.)



$$N = (-1)^s [1.m] \times 2^{E-127}$$

- E - 127 = e
- E = e + 127
- What is stored is E

Single Precision

- 23-bit mantissa
 - m ranges from 1 to $(2 - 2^{-23})$
- 8-bit exponent
 - E ranges from 1 to 254
 - 0 & 255 reserved to represent $-\infty$, $+\infty$, NaN
 - e ranges from -126 to +127
- Maximum representable number
 - $2 \times 2^{127} = 2^{128}$

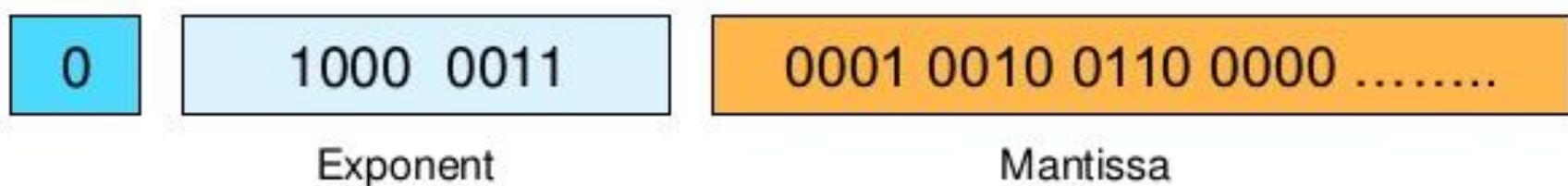
Example – Single Precision

- IEEE Single Precision Representation of 17.15_{10}
- Find 17.15_{10} in binary

$$\begin{aligned}17.15_{10} &= 10001.0010011_2 \\&= 1.00010010011 \times 2^4\end{aligned}$$

$$E = e + 127$$

$$E = 4 + 127 = 131$$



Floating-Point Numbers

- Floating-point number systems circumvent the limitation of having a constant number of integer and fractional bits
 - They allow the representation of very large and very small numbers
- The binary point **floats** to the right of the most significant 1
 - Similar to decimal scientific notation
 - For example, write 273_{10} in scientific notation:
 - Move the decimal point to the right of the most significant digit and increase the exponent:
$$273 = 2.73 \times 10^2$$
- In general, a number is written in scientific notation as:
$$\pm M \times B^E$$

Where,

- M = mantissa
- B = base
- E = exponent
- In the example, M = 2.73, B = 10, and E = 2 (that is, $+2.73 \times 10^2$)

Floating-Point Numbers

- Floating-point number representation using 32 bits
 - 1 sign bit
 - 8 exponent bits
 - 23 bits for the mantissa.



- The following slides show **three** versions of floating-point representation with 228_{10} using a 32-bit
 - The final version is called the **IEEE 754 floating-point standard**

Floating-Point Representation #1

- First, convert the decimal number to binary
 - $228_{10} = 11100100_2 = 1.11001 \times 2^7$
- Next, fill in each field in the 32-bit:
 - The sign bit (1 bit) is positive, so 0
 - The exponent (8 bits) is 7 (111)
 - The mantissa (23 bits) is 1.11001

1 bit	8 bits	23 bits
Sign	Exponent	Mantissa
0	00000111	11 1001 0000 0000 0000 0000

Floating-Point Representation #2

- You may have noticed that the first bit of the mantissa is always 1, since the binary point floats to the right of the most significant 1
 - Example: $228_{10} = 11100100_2 = 1.11001 \times 2^7$
- Thus, storing the most significant 1 (also called the implicit leading 1) is redundant information
- We can store just the **fraction parts** in the 23-bit field
 - Now, the leading 1 is implied

1 bit	8 bits	23 bits
Sign	Exponent	Fraction
0	0 0 0 0 0 1 1 1	110 0100 0000 0000 0000 0000

Floating-Point Representation #3

- The exponent needs to represent both positive and negative
- The final change is to use a **biased exponent**
 - The IEEE 754 standard uses a bias of 127
 - Biased exponent = bias + exponent
 - For example, an exponent of 7 is stored as $127 + 7 = 134 = 10000110_2$
- Thus , 228_{10} using the IEEE 754 32-bit floating-point standard is

Sign	Biased Exponent	Fraction
0	10000110	110 0100 0000 0000 0000 0000

Example

- Represent -58_{10} using the IEEE 754 floating-point standard
 - First, convert the decimal number to binary
 - $58_{10} = 111010_2 = 1.1101 \times 2^5$
 - Next, fill in each field in the 32-bit number
 - The sign bit is negative (1)
 - The 8 exponent bits are $(127 + 5) = 132 = 10000100_{(2)}$
 - The remaining 23 bits are the fraction bits: $11010000..000_{(2)}$

1 bit	8 bits	23 bits
Sign	Exponent	Fraction
1	10000100	110 1000 0000 0000 0000 0000

- It is $0xC2680000$ in the hexadecimal form
 - Check this out with the result of the sample program in the slide# 3

Floating-Point Number Precision

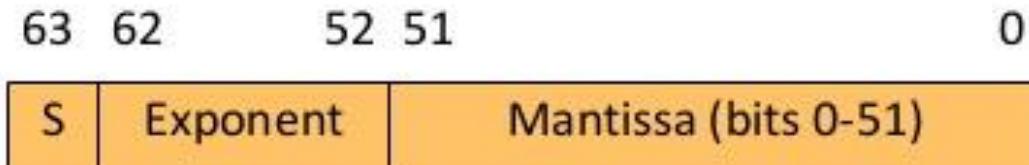
- The IEEE 754 standard also defines 64-bit double-precision that provides greater precision and greater range
 - **Single-Precision** (use the `float` declaration in C language)
 - 32-bit notation
 - 1 sign bit, 8 exponent bits, 23 fraction bits
 - bias = 127
 - It spans a range from $\pm 1.175494 \times 10^{-38}$ to $\pm 3.402824 \times 10^{38}$
 - **Double-Precision** (use the `double` declaration in C language)
 - 64-bit notation
 - 1 sign bit, 11 exponent bits, 52 fraction bits
 - bias = 1023
 - It spans a range from $\pm 2.22507385850720 \times 10^{-308}$ to $\pm 1.79769313486232 \times 10^{308}$
- Most general purpose processors (including Intel and AMD processors) provide hardware support for double-precision floating-point numbers and operations

Double Precision Example

- Represent -58_{10} using the IEEE 754 double precision
 - First, convert the decimal number to binary
 - $58_{10} = 111010_2 = 1.1101 \times 2^5$
 - Next, fill in each field in the 64-bit number
 - The sign bit is negative (1)
 - The 11 exponent bits are $(1023 + 5) = 1028 = 10000000100_{(2)}$
 - The remaining 52 bits are the fraction bits: $11010000..000_{(2)}$
 - It is $0xC04D0000_00000000$ in the hexadecimal form
 - Check this out with the result of the sample program in the slide# 4

IEEE Floating Point Standard (FPS)

- 2 standards
 - 1. Single precision
 - 32-bits
 - 23-bit mantissa
 - 8-bit exponent
 - 1-bit sign
 - 2. Double precision
 - 64-bits
 - 52-bit mantissa
 - 11-bit exponent
 - 1-bit sign



Single Precision

- 23-bit mantissa
 - m ranges from 1 to $(2 - 2^{-23})$
- 8-bit exponent
 - E ranges from 1 to 254
 - 0 & 255 reserved to represent $-\infty$, $+\infty$, NaN
 - e ranges from -126 to +127
- Maximum representable number
 - $2 \times 2^{127} = 2^{128}$

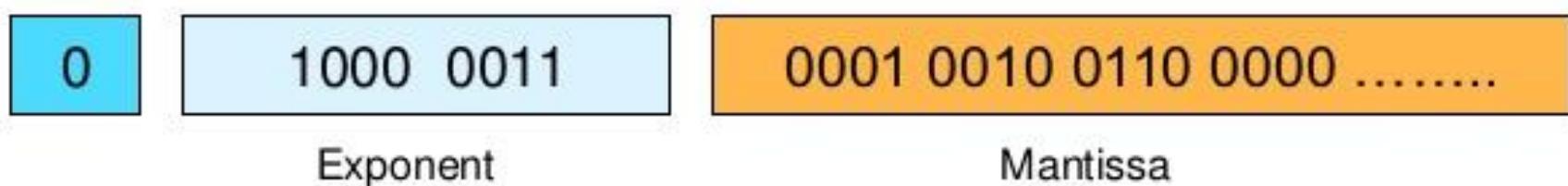
Example – Single Precision

- IEEE Single Precision Representation of 17.15_{10}
- Find 17.15_{10} in binary

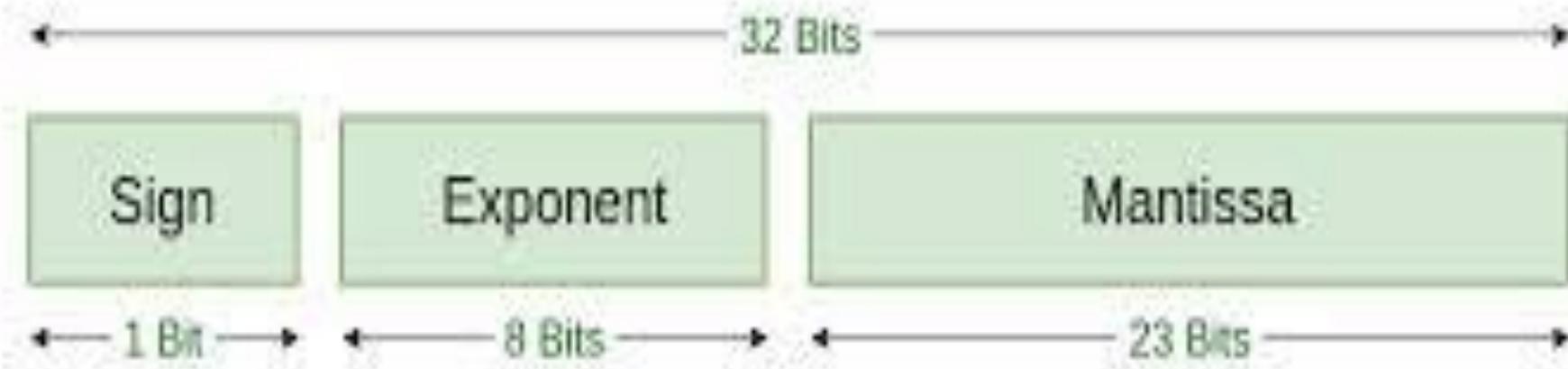
$$\begin{aligned}17.15_{10} &= 10001.0010011_2 \\&= 1.00010010011 \times 2^4\end{aligned}$$

$$E = e + 127$$

$$E = 4 + 127 = 131$$



IEEE standard Single precision



Single Precision
IEEE 754 Floating-Point Standard

Single precision

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then $S=0$ if not $S=1$ i.e.-Ve
4. Bias single precision format $E'=E+127$
5. Represent E' into binary number.
6. Put S, E and M values in single precision format.

Single precision

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 0101010010000000000000000

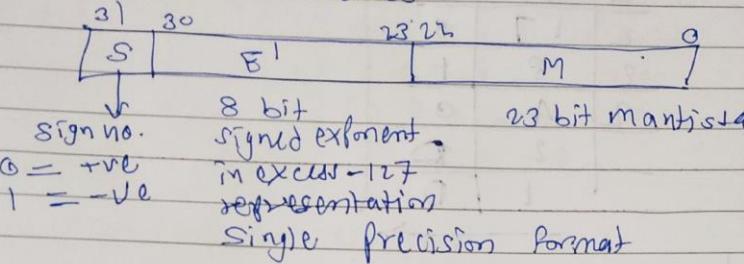
Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias single precision format $E'=E+127$
5. Represent E' into binary number.
6. Put S,E and M values in single precision format.

DATE: / /

IEEE Standard for Floating point Numbers
IEEE - Institute of Electrical and Electronic Engineers.

IEEE 754 standards



10.125 Represent in single precision formats.
Convert given number into binary

$$\begin{array}{r} 2 \mid 10 \\ 2 \mid 5 \quad | 0 \\ 2 \mid 2 \quad | 1 \\ 1 \quad | 0 \end{array} \quad \uparrow \quad \left. \right\} \text{Integer Part}$$

$$\begin{aligned} 0.125 \times 2 &= 0.25 & 0 & \downarrow & \left. \right\} \text{Fraction Part.} \\ 0.25 \times 2 &= 0.5 & 0 & \downarrow & \\ 0.5 \times 2 &= 1.0 & 1 & \downarrow & \\ 0 & & & & \end{aligned}$$

$$(10.125)_{10} = 1010.001$$

Normalize the number.

$$1010.001 = 1.01001 \times 2^3$$

S	E'	M
---	----	---

S = 0 +ve number

$$E = 3 \quad M = 01001$$

Bias for single precision format is = 127

$$E' = E + 127$$

$$E' = 3 + 127$$

$$E' = 13710$$

Steps for solving single precision problems

1. Convert given decimal number into binary (integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e. -Ve
4. Bias single precision format $E' = E + 127$
5. Represent E' into binary number.
6. Put S, E and M values in single precision format.

③ Represent S^1 into Binary $S^1 = 137$

2	137	1
2	68	1
2	34	0
2	17	0
2	8	1
2	4	0
2	2	0
1	0	.



$$E^1 = 1000100_2$$

S

E^1

M

0	1000 100	010001 --
---	----------	-----------

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then $S=0$ if not $S=1$ i.e.-Ve
4. Bias single precision format $E'=E+127$
5. Represent E' into binary number.
6. Put S, E and M values in single precision format.

Q) Represent 0.0625_{10} in single precision format?
 Ans. ① convert given decimal number in binary.

integer part = 0

fraction part

$$\begin{array}{r}
 0.0625 \times 2 = 0.125 \\
 0.125 \times 2 = 0.250 \\
 0.250 \times 2 = 0.5 \\
 0.5 \times 2 = 1.0
 \end{array}
 \quad \begin{array}{c}
 0 \\
 0 \\
 0 \\
 1
 \end{array}$$

$$0.0625_{10} = 0.0001_2$$

② Normalize the number

$$0.0001 = 1.0 \times 2^{-4}$$

③ Single precision representation

$$S=0 \quad E=-4 \quad m=0001$$

$$\begin{aligned}
 E' &= E + 127 \\
 &= -4 + 127 \\
 &= 123_{10}
 \end{aligned}$$

④ $E' = 123_{10}$ convert into binary

2	123	
2	6 ↑	1
2	3 0	1
2	1 5	0
2	7	1
2	3	1
	1	1

$$E' = 1111011_2$$

∴ numbers in single precision format

S	E'	M
0	0111011	00010-----

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias single precision format $E'=E=127$
5. Represent E' into binary number.
6. Put S,E and M values in single precision format.

- Q2 Represent 15.2_{10} in single precision format.
- Ans. ① convert given number in single precision format

integer part

$$\begin{array}{r} 2 \mid 15 \\ 2 \mid 7 \\ 2 \mid 3 \\ 1 \end{array} \quad \rightarrow$$

fraction part

$$\begin{array}{l} 0.2 \times 2 = 0.4 \\ 0.4 \times 2 = 0.8 \\ 0.8 \times 2 = 1.6 \\ 0.6 \times 2 = 1.2 \\ 0.2 \times 2 = 0.4 \end{array} \quad \begin{array}{c} 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{array}$$

$$15.2_{10} = 1111.00110_2$$

② Normalize the number.

$$1111.00110 = 1.11100110 \times 2^3$$

③ Represent in single precision format:

$$S=0 \quad E=3 \quad m=11100110$$

$$E' = E + 127 \\ = 3 + 127 = 130$$

④ convert E' into (130_{10}) into binary

$$\begin{array}{r} 2 \mid 130 \\ 2 \mid 65 \quad 0 \\ 2 \mid 32 \quad 1 \\ 2 \mid 16 \quad 0 \\ 2 \mid 8 \quad 0 \\ 2 \mid 4 \quad 0 \\ 2 \mid 2 \quad 0 \\ 1 \end{array} \quad \begin{array}{c} 130_{10} = (10000010)_2 \\ S \quad E' \quad . \quad M \end{array}$$

Steps for solving single precision problems

1. Convert given decimal number into binary (integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then $S=0$ if not $S=1$ i.e. -Ve
4. Bias single precision format $E'=E+127$
5. Represent E' into binary number.
6. Put S, E and M values in single precision format.

- Q) Represent 28.5 into single precision format.
 Ans. ① convert given number into binary.

Integer part

2	28	
2	14	0
2	7	0
2	3	1
1		1

Fraction part

$$28.5 = \underbrace{11100}_{\text{M}}.\underbrace{10}_{\text{E}}$$

$0.5 \times 2 = 1.0$
 $0 \times 2 = 0$

↓
Normalization
 1.110010×2^7

- ② Single precision format

$$S = 1 \text{ because number is -ve } E = 4$$

$$E' = E + 127$$

$$= 4 + 127$$

$$E' = 131_{10}$$

- ③ Convert into binary

2	131	
2	65	1
2	32	1
2	16	0
2	8	0
2	4	0
2	2	0
1		0

$$E' = 1000001_2$$

S	E	M
1	10000011	110010...

Steps for solving single precision problems

1. Convert given decimal number into binary (integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e. -Ve
4. Bias single precision format $E' = E + 127$
5. Represent E' into binary number.
6. Put S, E and M values in single precision format.

(a) 100.125

Convert it into binary

$$(100)_{10} = (1100100)_2$$

Fraction part

$$0.125 \times 2 = 0.250$$

$$0.25 \times 2 = 0.50$$

$$0.50 \times 2 = 1.00$$

$$\therefore (100.125)_{10} = (1100100.001)_2$$

To convert above no in single precision format we will convert

$$N \times 2^{E-127}$$

$$1100100.001 = 1.100100001 \times 2^6$$

$$= 1.100100001 \times 2^{133-127}$$

0	1000101	1001 100000 100000
sign	E	significant (B)

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias single precision format $E' = E + 127$
5. Represent E' into binary number.
6. Put S,E and M values in single precision format.

(b) 42.625

$$(42)_{10} = (101010)_2$$

$$0.625 \times 2 = 1.350$$

$$0.350 \times 2 = 0.70$$

$$0.70 \times 2 = 1.40$$

$$0.40 \times 2 = 0.80$$

$$0.80 \times 2 = 1.60$$

$$0.60 \times 2 = 1.20$$

$$0.20 \times 2 = 0.40$$

$$\therefore (0.625)_{10} = 101011001100110$$

$$42.625 = 101010.10101100110$$

Convert it into $1.N \times 2^{E-127}$

$$1.0101010101100110 \dots \times 2^5$$

$$= 1.0101010101100110 \times 2^{132-127}$$

0	10000100	01010101011011001100110	mantissa
sign bit	Exponent		

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias single precision format $E'=E+127$
5. Represent E' into binary number.
6. Put S,E and M values in single precision format.

(c) 17.25

$$(17)_{10} = (10001)_2$$

$$0.125 \times 2 = 0.50$$

$$0.50 \times 2 = 1.00$$

$$\therefore (17.125)_{10} = (10001.01)_2$$

Convert it into the form $(1.N) \times 2^{E-127}$

$$10001.01 = 1.000101 \times 2^4$$

$$= 1.000101 \times 2^{131-127}$$

0	10000011	0001 0100 0000 0000 0000
sign	Exponent	significand (M)

Steps for solving single precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias single precision format E'=E+127
5. Represent E' into binary number.
6. Put S,E and M values in single precision format.

Represent $(178.1875)_{10}$ in single and double precision floating point format.

Step I : Convert decimal number in binary format integer part first

$$\begin{array}{r} 11 \rightarrow B \\ 16 \overline{) 178} \\ 16 \\ \hline 18 \\ 16 \\ \hline 2 \end{array}$$

$$(178)_{10} = (B2)_{16} = (\underbrace{1\ 0\ 1\ 1}_{B} \underbrace{0\ 0\ 1\ 0}_2)_2$$

Fractional Part :

$$0.1875 \times 2 = 0.3750 \quad 0$$

$$0.3750 \times 2 = 0.750 \quad 0$$

$$0.750 \times 2 = 1.51$$

$$0.5 \times 2 = 1.01$$

$$0 \quad 0$$

$$= 0.0011$$

$$\therefore \text{Binary number} = 1011\ 0010 + 0.0011 = 1011\ 0010. 0011$$

IS = 1022 <= E <= 1023.

Example 1 : Represent 1259.125_{10} in single precision and double precision formats.

Solution : Step 1 : Convert decimal number in binary format.

$$\begin{array}{r} 78 \\ \hline 16) 1259 \\ 112 \\ \hline 0139 \end{array}$$

$$\begin{array}{r} 4 \\ \hline 16) 78 \\ 64 \\ \hline 14 = E \end{array}$$

$$\begin{array}{r} 128 \\ \hline 011 = M \end{array}$$

$$= 4EBH = \frac{100}{4} = \frac{1110}{E} \frac{1011}{B}$$

$$0.125 \times 2 = 0.25 \text{ O}$$

$$0.25 \times 2 = 0.5 \text{ O}$$

$$0.5 \times 2 = 1.0 \text{ 1}$$

O

$$= 0.001$$

$$\therefore \text{Binary number} = 10011101011 + 0.001$$

$$= 10011101011.001$$

Step 2 : Normalize the number

$$10011101011.001 = 1.0011101011001 \times 2^{10}$$

Now we will see the representation of the numbers in single precision and double precision formats.

Step 3 : Single precision representation for a given number S = 0, E = 10 and M = 00111 01011 001

Bias for single precision format is = 127

$$\begin{aligned} E' &= E + 127 = 10 + 127 = 137_{10} \\ &= 10001001_2 \end{aligned}$$

\therefore Number in double precision format is given as

Sign $\underbrace{1000 \ 1001}_{\text{Exponent}}$ $\underbrace{0011101011001 \dots 0}_{\text{Mantissa (23 bits)}}$

Step 4 : Double precision representation for a given number

S = 0, E = 10, and M = 0011101011001

Bias for double precision format is = 1023

$$\begin{aligned} E' &= E + 1023 = 10 + 1023 = 1033_{10} \\ &= 10000001001_2 \end{aligned}$$

\therefore Number in double precision format is given as

Sign $\underbrace{1000 \ 1001}_{\text{Exponent}}$ $\underbrace{0011101011001 \dots 0}_{\text{Mantissa (52 bits)}}$

Example 2 : Represent -307.1875_{10} in single precision and double precision formats.

Solution : **Step 1 :** Convert decimal number in binary format integer part.

Integer part :

$$\begin{array}{r} 19 \\ 16 \overline{) 307} \\ 16 \\ \hline 147 \\ 144 \\ \hline 3 = 3 \end{array}$$

$$\begin{array}{r} 1 \\ 16 \overline{) 19} \\ 16 \\ \hline 3 \end{array}$$

$$= 133H = \frac{1}{1} = \frac{0011}{3} \quad \frac{0011}{3}$$

Fractional part :

$$0.1875 \times 2 = 0.3750 \quad 0$$

$$0.3750 \times 2 = 0.750 \quad 0$$

$$0.750 \times 2 = 1.5 \quad 1$$

$$0.5 \times 2 = 1.0 \quad 1$$

$$0$$

$$= 0.0011$$

$$\begin{aligned} \text{Binary number} &= -100110011 + 0.0011 \\ &= -100110011.0011 \end{aligned}$$

Step 2 : Normalize the number -100110011.0011
 $= -1.001100110011 \times 2^8$

Now we will see the representation of the numbers in single precision and double precision formats.

Step 3 : Single precision representation for a given number

$$S = 1, E = 8, \text{ and } M = 0011 0011 0011$$

Bias for single precision format is $= 127$

$$\begin{aligned} E' &= E + 127 = 8 + 127 = 135_{10} \\ &= 10000111_2 \end{aligned}$$

Number in single precision format is given as

Sign Exponent Mantissa (23 bits)

$\underbrace{1000}_{\text{Sign}}$ $\underbrace{0111}_{\text{Exponent}}$

$\underbrace{00110011001100 \dots 0}_{\text{Mantissa (23 bits)}}$

Step 2 : Normalize the number

$$1011\ 0010.0011 = 1.01100100011 \times 2^7$$

Step 3 : Representation in signal precision for given number

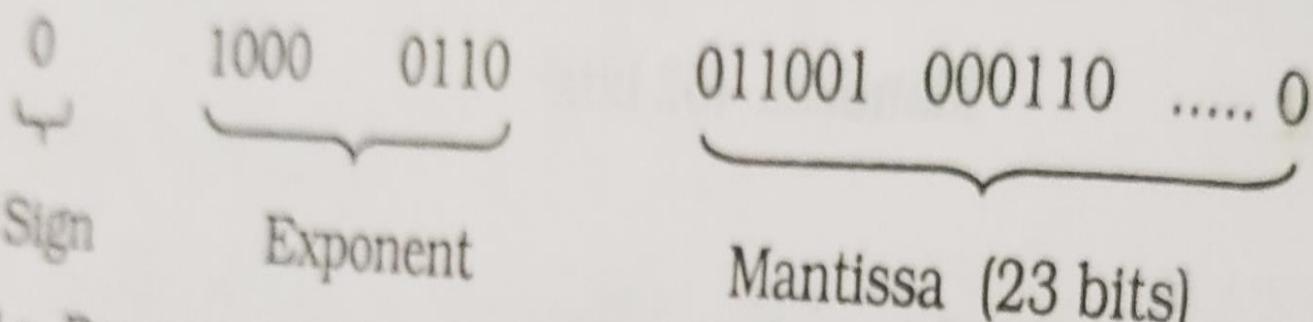
$$S = 0, E = 7 \text{ and } M = 011\ 001\ 000\ 11$$

Bias for single precision format is = 127

$$E' = E + 127 = 7 + 127 = 134_{10}$$

$$= 10000110_2$$

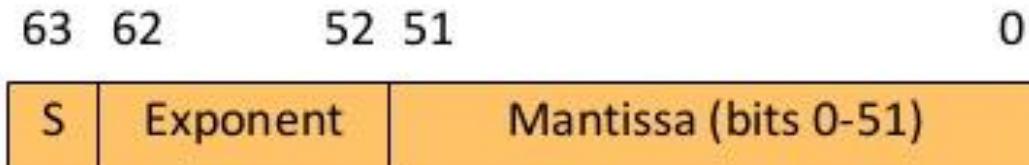
∴ Number in single precision format is given as



Step 4 : Representation in double

IEEE Floating Point Standard (FPS)

- 2 standards
 - 1. Single precision
 - 32-bits
 - 23-bit mantissa
 - 8-bit exponent
 - 1-bit sign
 - 2. Double precision
 - 64-bits
 - 52-bit mantissa
 - 11-bit exponent
 - 1-bit sign



Representation of Floating Point Numbers in Single Precision IEEE 754 Standard

$$\text{Value} = N = (-1)^{\text{S}} \times 2^{\text{E}-127} \times (1.\text{M})$$

$-128 < E < 255$
Actual exponent is:
 $\text{e} = E + 127$



exponent: excess 127
binary integer added

mantissa: sign + magnitude, normalized
binary significand with
a hidden integer bit: 1.M

Example: 0 = 0 00000000 0...0 -1.5 = 1 01111110...0

Magnitude of numbers that can be represented is in the range:

$$2^{-126} (1.0) \text{ to } 2^{127} (2 \cdot 2^{-23})$$

Which is approximately:

$$1.8 \times 10^{-38} \text{ to } 3.40 \times 10^{38}$$

Representation of Floating Point Numbers in Double Precision IEEE 754 Standard

$$\text{Value} = N = (-1)^S \times 2^{E-1023} \times (1.M)$$



$$0 < E < 2047$$

Actual exponent is:
 $e = E - 1023$

exponent:
excess 1023
binary integer
added

mantissa:
sign + magnitude, normalized
binary significand with
a hidden integer bit: 1.M

Example: $0 = 0\ 0000000000\ 0\dots 0$ $-1.5 = 1\ 0111111111\ 10\dots 0$

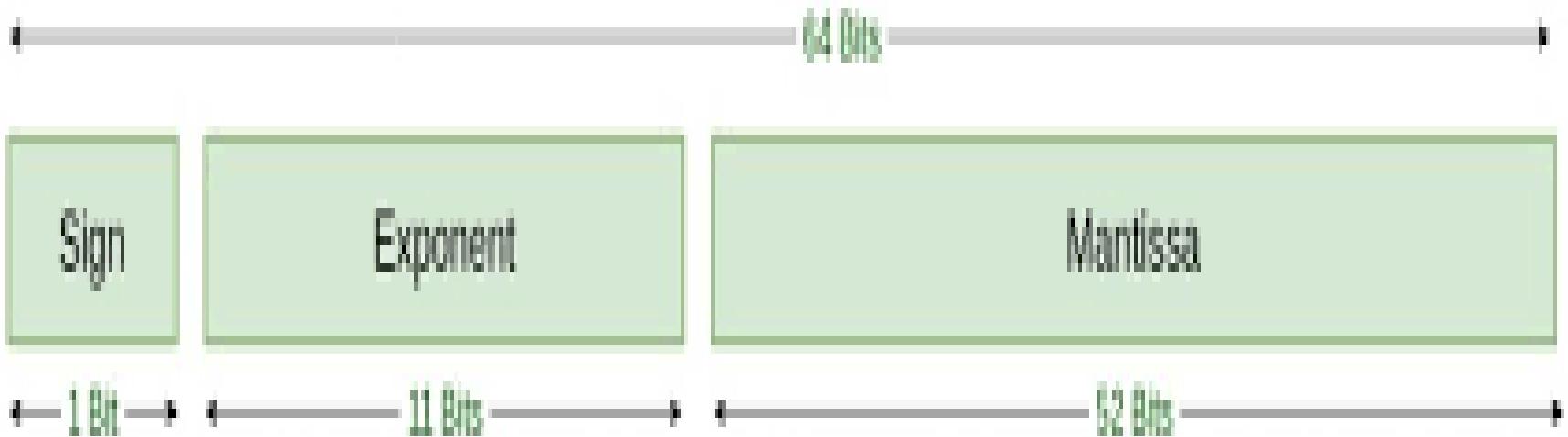
Magnitude of numbers that
can be represented is in the range:

$$2^{-1022} (1.0) \text{ to } 2^{1023} (2 \cdot 2^{52})$$

Which is approximately:

$$2.23 \times 10^{-308} \text{ to } 1.8 \times 10^{308}$$

Double Precision-IEEE754 format



Double Precision
IEEE 754 Floating-Point Standard

Double Precision

Steps for solving double precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then $S=0$ if not $S=1$ i.e.-Ve
4. Bias double precision format $E'=E+1023$
5. Represent E' into binary number.
6. Put S, E and M values in double precision format.

① Represent the following number in double precision format.

$$100.125_{10}$$

Step 1 convert given number into binary

2	100	
2	50	0
2	25	0
2	12	1
2	6	0
2	3	0
2	1	1

integer part.

$$0.125 \times 2 = 0.250$$

$$0.250 \times 2 = 0.500$$

$$0.500 \times 2 = 1.00$$

0 } | { fraction part.
0 } | { 1 |

Step 2 Normalized the number.

$$\begin{aligned}100.125_{10} &= 1100100 \cdot 001 \\&= 1.100100001 \times 2^6\end{aligned}$$

Step 3 S = 0 E = 6 M = 10010001

$$\begin{aligned}E' &= E + 1023 \\&= 6 + 1023 = 1029_{10}\end{aligned}$$

Step 4 convert E' into binary

-2	1029		1029 ₁₀ = 11001000101 ₂
2	514	1	S
2	257	0	E
2	128	1	M
2	64	0	1100000000010110010001 ---
2	32	0	16 bit. 11 bits 52 bits
2	16	0	
2	8	0	
2	4	0	

Steps for solving double precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias double precision format E'=E+1023
5. Represent E' into binary number.
6. Put S,E and M values in double precision format.

② 42.625_{10} double precision.

Step 1

2	42		
2	21	0	
2	10	1	↑
2	5	0	
2	2	1	
1	0		

↓ ↓

Integer Part Fraction part

$0.625 \times 2 = 1.250$ 1
 $0.250 \times 2 = 0.500$ 0
 $0.500 \times 2 = 1.000$ 1
 $0.000 \times 2 = 0.000$ 0
 $0.000 \times 2 = 0.000$ 0
 $0.000 \times 2 = 0.000$ 0

Step 2

$$42.625_{10} = 101010 \cdot 1010100101110_2 \times 2^5$$

Step 3

$$S = 0 \quad E = 5 \quad M = 01010101011001100110$$

$$E' = E + 1023$$

$$= 5 + 1023 = 1028_{10}$$

Step 4.

2	1028	↓
2	514	0
2	257	0
2	128	1
2	64	0
2	32	0
2	16	0
2	8	0
2	4	0
2	2	0
1	1	0

$1028_{10} \leq 10000000100_2$

Step 5

S	E'	M
0	10000000100	0101010101101100110...

1 bit 11 bits 52 bits

Steps for solving double precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias double precision format $E'=E+1023$
5. Represent E' into binary number.
6. Put S,E and M values in double precision format.

③ 12.5 solve using double precision.

841	2	12	$0.5 \times 2 = 1.0$	↓
	2	6	$0 \times 2 = 0$	↓
	2	3	Fraction part	
	1	1		
			Integer Part.	

Step 2 Normalization.

$$12.5_{10} = 1100.10_2$$

$$= 1.10010 \times 2^3$$

8413 $s = 0$ $E = 3$ $M = 10010$

8414 $E' = E + 1023$

$$= 3 + 1023 = 1026_{10}$$

8415

2	1026	$\overline{-}$
2	513	0
2	206	1
2	103	0
2	51	1
2	25	1
2	12	1
2	6	0
2	3	0
1	1	

$$1026_{10} = 1100111010_2$$

S	E	M
0	1100111010	10010 -----
1 bit	11 bits	52 bits

Steps for solving double precision problems

1. Convert given decimal number into binary(integer and fraction part).
2. Normalize the binary number and find values of Sign, Exponent, Mantissa.
3. If given number is +ve then S=0 if not S=1 i.e.-Ve
4. Bias double precision format $E' = E + 1023$
5. Represent E' into binary number.
6. Put S,E and M values in double precision format.

IS = 1022 <= E <= 1023.

Example 1 : Represent 1259.125_{10} in single precision and double precision formats.

Solution : Step 1 : Convert decimal number in binary format.

$$\begin{array}{r} 78 \\ 16 \overline{) 1259} \\ 112 \\ \hline 0139 \end{array}$$

$$\begin{array}{r} 4 \\ 16 \overline{) 78} \\ 64 \\ \hline 14 = E \end{array}$$

$$\begin{array}{r} 128 \\ \hline 011 = M \end{array}$$

$$= 4EBH = \frac{100}{4} = \frac{1110}{E} \frac{1011}{B}$$

$$0.125 \times 2 = 0.25 \text{ O}$$

$$0.25 \times 2 = 0.5 \text{ O}$$

$$0.5 \times 2 = 1.0 \text{ 1}$$

O

$$= 0.001$$

$$\therefore \text{Binary number} = 10011101011 + 0.001$$

$$= 10011101011.001$$

Step 2 : Normalize the number

$$10011101011.001 = 1.0011101011001 \times 2^{10}$$

Now we will see the representation of the numbers in single precision and double precision formats.

Step 3 : Single precision representation for a given number S = 0, E = 10 and M = 00111 01011 001

Bias for single precision format is = 127

$$\begin{aligned} E' &= E + 127 = 10 + 127 = 137_{10} \\ &= 10001001_2 \end{aligned}$$

\therefore Number in double precision format is given as

Sign Exponent Mantissa (23 bits)

Step 4 : Double precision representation for a given number

S = 0, E = 10, and M = 0011101011001

Bias for double precision format is = 1023

$$\begin{aligned} E' &= E + 1023 = 10 + 1023 = 1033_{10} \\ &= 10000001001_2 \end{aligned}$$

\therefore Number in double precision format is given as

Sign Exponent Mantissa (52 bits)