

UNIT II

①

Arithmetic Operations:-

ALU - Addition and Subtraction - Multiplication
- division - floating point operation - Subword
parallelism.

① ALU (Arithmetic & Logic Unit)

1.1 Arithmetic & Logic Unit:-

A basic operation in ~~all~~ digital computers is the addition or subtraction of two numbers. Arithmetic operations occur at the machine instruction level. They are implemented along with basic logic functions such as AND, OR, NOT and Exclusive OR (XOR), in the Arithmetic and Logic Unit (ALU) subsystem of the processor.

The time needed to perform an addition operation affects the processors performance. Multiply and divide operation which require more complex circuitry than either addition or subtraction operation also affect the performance.

1.2) Signed and Unsigned Numbers:-

Computer represent numbers in binary. The hexadecimal notation is sometimes used to represent binary number as it is more concise.

Binary numbers has two possible digits or bit:

0 and 1. So they are considered base 2 number.

(base 10 for decimal). Generalizing the point in

any number base the value of i th digit d is

$$\boxed{d \times \text{base}^i}$$

Where i starts at 0 and increases from right to left.

In real world of Mathematics computers must represent both positive and negative binary numbers. For eg: even when dealing with positive arguments, mathematical operations may produce a negative result. (ie) When we subtract larger number from smaller number

eg: $124 - 237 = -113.$

Thus a consistent method for representing negative numbers in binary computer arithmetic operation is needed.

There are various approach but they all involve ^② using one of the digit of the binary number to represent the sign of the number. There are three methods, They are the ① Sign and Magnitude representation, ② One's Complement Representation, ③ Two's Complement Representation.

1.2.1) Sign and Magnitude Representation:-

In sign and Magnitude representation, the left bit is used to reserved for sign bit (0 indicates positive number & 1 indicates negative numbers).

Sign bit	All bit Right are the number Magnitude
----------	--

↓
if 0 +ve number
if 1 -ve number

Advantages of Sign - Magnitude:-

- ① It is very simple to implement
- ② Useful for floating point representation.

Disadvantages of sign - Magnitude:-

- ① Sign bit is independent of Magnitude; can be both +0, -0 difficult to represent.

1.2) Ones Complement Representation:

The negative number of the same Magnitude as any given positive number is its ones complement. If $m = 01001100$ then m complement is 10110011 . The most significant bit is the sign and is 0 for positive binary number and 1 for negative number.

Drawbacks:

The drawback in ones complement method has two zeros (+0 & -0) i.e.

$$m \quad 00000000 = +0$$

$$11111111 = -0$$

it is hard to implement in hardware.

1.2.3) Two's Complement Representation:

Since there is drawbacks in Sign & Magnitude representation and ones complement, another approach called the 2's complement has become a standard for representing the sign of a fixed point binary number in computer circuits. The 2's complement representation is widely used in computers and it is a signed representation as it can represent both positive and

negative numbers.

3

In general, Most Significant bit (MSB or S) is 0 indicate a positive number and it is represented by its binary equivalent directly. If S=1 indicate a negative number or represent 2's Complement Value. The 2's complement of a negative number can be obtained by adding one to the one's complement.

eg: The two's complement of -13 can be obtained as follow.

(*) Step 1 :- Obtain the binary of +13 [8 bit representation]

(*) Step 2: Take one's complement
1's complement is: 1111 0010

(*) Step 3: Add 1 to 1's complement.

$$\begin{array}{r} 1111\ 0010 \\ +1 \\ \hline -13 = 1111\ 0011_2 \end{array}$$

$$\begin{array}{r} 2 \overline{) 13} \\ 2 \overline{) 6} - 1 \\ 2 \overline{) 3} - 0 \\ 1 - \phi \end{array}$$

$(13)_{10} = 1101_2$

Which is the Value of -13!

The 2's complement representation of a negative number always has '1' as its MSB.

In a n -bit representation, if we add $2n$ -bit numbers result is $n+1$ bits the left most bit is discarded. Since in computer there are no extra columns. If there are only 8 bit the 9th bit if resulted will be discarded.

eg: Consider 4 bit Addition.

$$\begin{array}{r} 12_{10} + 14_{10} \\ 1100_{10} \\ + 1110_2 \\ \hline 111010 \\ \swarrow \text{discarded} \end{array}$$

because it can hold only 4 bit.

1.3 Addition and Subtraction:

Addition is carried out similar to decimal addition (ie) start adding bit by bit from right to left, carry generated during addition, then it is taken to next digit to the left (higher order bit).

Subtraction uses addition (2's complement)
The appropriate operand is negated before being added.

3.1) Boolean Addition:-

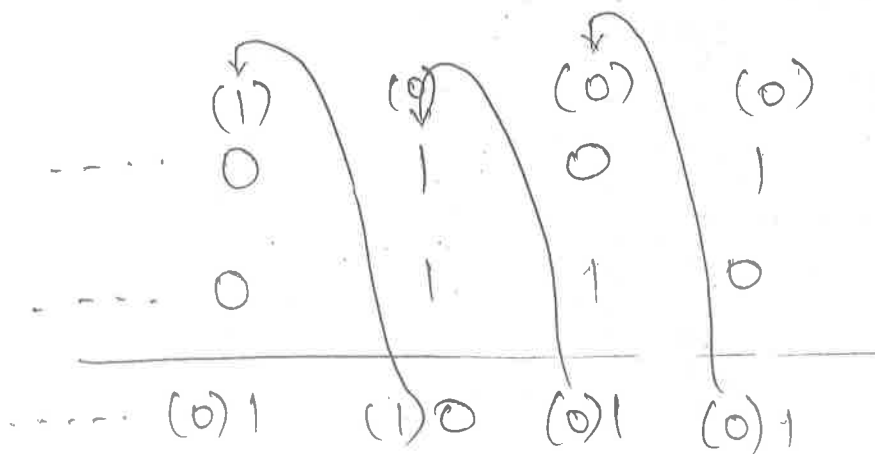
4

When adding two numbers, if the sum of the digit in a given position equals or exceeds, then a carry is propagated. For eg: if two ones are added sum is 10_2 thus record a 0 for the sum and propagate a carry value 1 into next higher significant bit.

eg:-

$$5_{10} + 6_{10}$$

$$\begin{array}{r}
 \begin{array}{cccccccc}
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0101_2 \text{ (5)}_{10} \\
 (+) & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0110_2 \text{ (6)}_{10} \\
 \hline
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 1011_2 \text{ (11)}_{10}
 \end{array}
 \end{array}$$



Carry propagation

1.3.2) Boolean Subtraction:-

There are two ways of performing boolean subtraction:

1) Using normal subtraction procedure.

2) negate the subtrahend ($a - b$ is $a + \text{neg } b$ ~~subtrahend~~)

then perform addition.

eg: $5_{10} - 6_{10}$

\Downarrow

$5_{10} + (-6_{10})$

$\rightarrow -b$ can be obtained by 2's complement method.

eg:-

i) $8_{10} - 5_{10}$

$$\begin{array}{r}
 \begin{array}{cccccccc}
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0100 \text{ (} 8_{10} \text{)} \\
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0101 \text{ (} 5_{10} \text{)} \\
 \hline
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0011 \text{ (} 3_{10} \text{)} \\
 \hline
 \end{array}
 \end{array}$$

(or) by using 2nd Method. (ie 2's complement)

ii) $8_{10} - 5_{10} \Rightarrow 8_{10} + (-5_{10})$

find $-5_{10} \Rightarrow$

ones complement \Rightarrow

$$\begin{array}{r}
 \begin{array}{cccccccc}
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0101 \text{ (+5)} \\
 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1010 \\
 & & & & & & & +1 \\
 \hline
 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1011 \text{ (-5)}
 \end{array}
 \end{array}$$

negative numbers.

3

In general, Most Significant bit (MSB or S) is 0 indicate a positive number and it is represented by its binary equivalent directly. If S=1 indicate a negative number or represent 2's Complement Value. The 2's complement of a negative number can be obtained by adding one to the one's complement.

eg: The two's complement of -13 can be obtained as follow.

(*) Step 1 :- Obtain the binary of +13 [8 bit representation]

(*) Step 2: Take one's complement
1's complement is: 1111 0010

(*) Step 3: Add 1 to 1's complement:

$$\begin{array}{r} 1111\ 0010 \\ +1 \\ \hline -13 = 1111\ 0011_2 \end{array}$$

$$\begin{array}{r} 2 \overline{) 13} \\ 2 \overline{) 6} - 1 \\ 2 \overline{) 3} - 0 \\ 1 - \phi \end{array}$$

$(13_{10}) = 1101_2$

Which is the Value of -13

The 2's complement representation of a negative number always has '1' as its MSB.

In a n -bit representation, if we add $2n$ -bit numbers result is $n+1$ bits the left most bit is discarded. Since in computer there are no extra columns. If there are only 8 bit the 9th bit if resulted will be discarded.

eg: Consider 4 bit Addition.

$$\begin{array}{r} 12_{10} + 14_{10} \\ 1100_{10} \\ + 1110_2 \\ \hline 111010 \\ \swarrow \text{discarded} \end{array}$$

because it can hold only 4 bit.

1.3 Addition and Subtraction:

Addition is carried out similar to decimal addition (ie) start adding bit by bit from right to left, carry generated during addition, then it is taken to next digit to the left (higher order bit).

Subtraction uses addition (2's complement)
The appropriate operand is negated before being added.

now add 8_{10} and $(-5)_{10}$

(5)

$$\begin{array}{r}
 \begin{array}{cccccccccccc}
 (1) & & & & & & & & & & (1)(1)(1)(1) \\
 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 1000 & (8_{10}) \\
 (+) & 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1111 & 1011 & (-5_{10}) \\
 \hline
 1 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0000 & 0011 & (3_{10}) \\
 \text{discarded} & & & & & & & & & &
 \end{array}
 \end{array}$$

4) Overflow:

" Overflow occurs when there is insufficient bit in the binary number representation to hold the result of an arithmetic operation.

When a 32 bit number is added or subtracted a result produced may need 33 bit to store the result. This condition is called overflow.

eg: Consider 4 bit

$$\begin{array}{rcl}
 8_{10} & 1000 & \rightarrow 4 \text{ bit} \\
 9_{10} & (+) 1001 & \rightarrow 4 \text{ bit} \\
 \hline
 \end{array}$$

$$(1)0001 \rightarrow 5 \text{ bit (4th bit is needed)}$$

This is called overflow bit.

To detect or compensate for overflow one needs $n+1$ bits, if an n -bit number representation is employed.

1.4.1) Conditions where Overflow cannot occur:-

- ① In addition if the operands are of different sign. [eg: $(+5) + (-4)$] \Rightarrow $\begin{array}{r} 0101 \text{ (5)} \\ + 1010 \text{ (-4)} \\ \hline 1111 \end{array}$ (no overflow) -4
 \Rightarrow 0100
1001
1010
- ② In subtraction if the operands are of same sign. [eg: $(+5) - (+2)$ or $(-5) - (-2)$]

1.4.2) Conditions where overflow occurs:-

Operation	operand A	operand B	result Indicating Overflow
Addition $\left\{ \begin{array}{l} A+B \\ A+B \end{array} \right.$	≥ 0	≥ 0	< 0
	< 0	< 0	≥ 0
Subtraction $\left\{ \begin{array}{l} A-B \\ A-B \end{array} \right.$	≥ 0	< 0	< 0
	< 0	≥ 0	≥ 0

(ie) i) Overflow occurs when we adding two positive numbers and sum is a negative or adding two negative numbers and result is positive.

ii) When we subtract a negative number from a positive number I get a negative result or when we subtract a positive number from a negative number and get a positive result.

1.4.3) MIPS overflow Handling:

- It raises an exception when overflow occurs
 - 2's Complement arithmetic operation (add, addi, sub raise exception overflow).
 - Addu and addiu do not raise exception, since they are used for arithmetic operation on addresses.
-

2. Addition & Subtraction.

(7)

Addition and subtraction are the Arithmetic operations performed in the ALU (Arithmetic Logic Unit). Addition is performed by adding bit by bit digits from right to left, carry generated during the addition is taken to next digit to the left. Subtraction can be performed by subtracting the digits or adding the 2's complement of the subtrahend).

eg:-

$$5_{10} + 6_{10}$$

$$\begin{array}{r}
 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0101_2 \\
 +\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0110_2 \\
 \hline
 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1011\ (11)_{10}
 \end{array}$$

Carry propagation

$$\begin{array}{r}
 \overset{(1)}{0} \overset{(0)}{1} \overset{(0)}{0} \overset{(0)}{1} \\
 0\ 1\ 1\ 0 \\
 \hline
 (0)\ 1\ (1)\ 0\ (0)\ 1\ (0)\ 1 \Rightarrow (1011)_2
 \end{array}$$

eg: $8_{10} - 5_{10} = (3)_{10}$

0000	0000	0000	0000	0000	0000	0000	1010 ₂
0000	0000	1000	0000	0000	0000	0000	0101 ₂
<hr/>							
0000	0000	0000	0000	0000	0000	0000	0011 ₂
<hr/>							2
							↓
							3 ₁₀

There are some Adders, High speed Adders for performing the operation.

- ① Full Adder
- ② Ripple carry Adder
- ③ Carry Lookahead Adder — Fast Adders.

Full Adder:-

Full Adder is a circuit for performing 1 bit addition, if two bits x_i and y_i are added Sum S_i and C_{i+1} is generated. A full adder uses three input one is x_i , another is y_i third input is the C_i represents the carry-in to the i th stage.

The logic truth table for Full adder is given below.

Truth Table for Full Adder.

8

Input			output	
x_i	y_i	Carry-in c_i	Sum s_i	Carry out c_{i+1}
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

The logical expression:-

$$\text{Sum } s_i = \bar{x}_i \bar{y}_i c_i + \bar{x}_i y_i \bar{c}_i + x_i \bar{y}_i \bar{c}_i + \cancel{x_i y_i c_i}$$

$$= \bar{x}_i (\bar{y}_i c_i + y_i \bar{c}_i) + x_i (\bar{y}_i \bar{c}_i + y_i c_i)$$

$$= \bar{x}_i (y_i \oplus c_i) + x_i (\bar{y}_i \oplus \bar{c}_i) \quad \left(\because A \oplus B = \bar{A}B + A\bar{B} \right)$$

$$\boxed{\text{Sum } s_i = x_i \oplus y_i \oplus c_i}$$

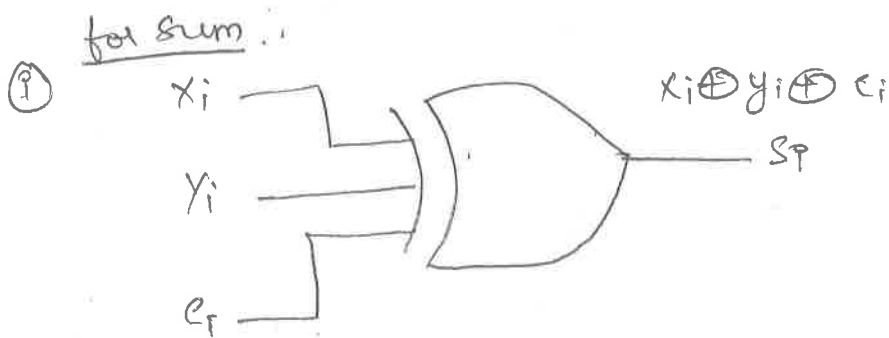
$$\begin{aligned} \text{Carry } c_{i+1} &= \bar{x}_i y_i c_i + x_i \bar{y}_i c_i + x_i y_i \bar{c}_i + x_i y_i c_i \\ &= y_i c_i (\bar{x}_i + x_i) + x_i c_i (\bar{y}_i + y_i) + x_i y_i (\bar{c}_i + c_i) \end{aligned}$$

$$\boxed{c_{i+1} = y_i c_i + x_i c_i + x_i y_i}$$

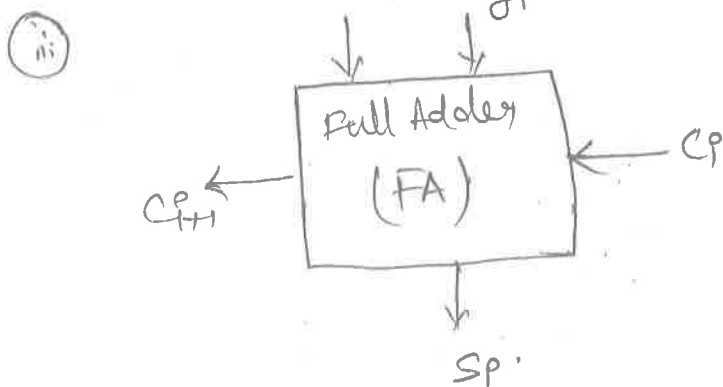
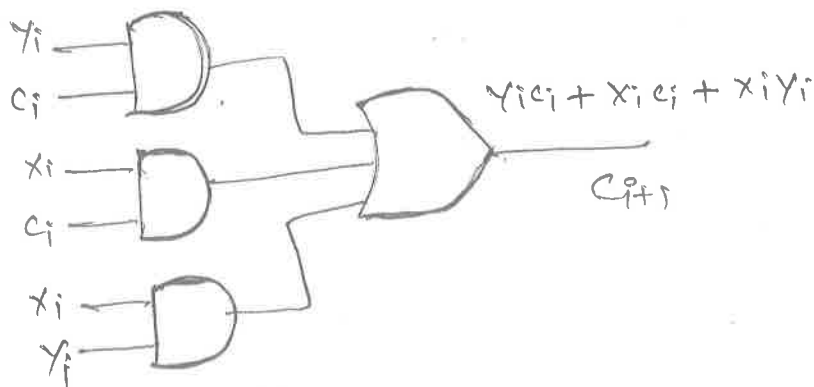
$$\because A + A = A$$

The logical Expression^{for sum} can be implemented using a 3-input XOR gate. And carry C_{i+1} can be implemented with two-level AND-OR logic circuit.

A conventional symbol for the complete circuit for a single stage of addition called Full Adder (FA) -



② for carry :



(9)

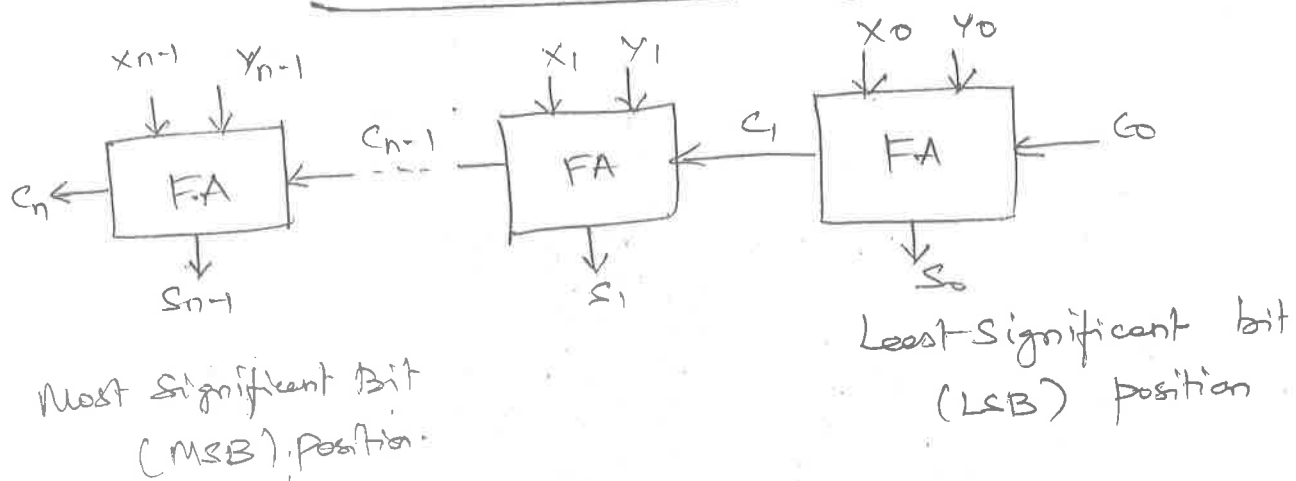
2) Ripple Carry Adder:-

A single full adder is capable of adding two one-bit numbers and an input carry.

In order to add binary numbers with more than one bit, additional full adders must be employed.

The n -bit parallel adder can be constructed using number of full adder circuit connected in parallel.

An n -bit ripple Carry Adder.



Each 1 bit adder stage supplies a carry bit to the stage on its left, hence carry signals propagate through the adder from Right to Left.

giving rise to the name "ripple carry adder".

In the worst case a carry signal can ripple through all n stage of adder. The input carry is normally set to 0.

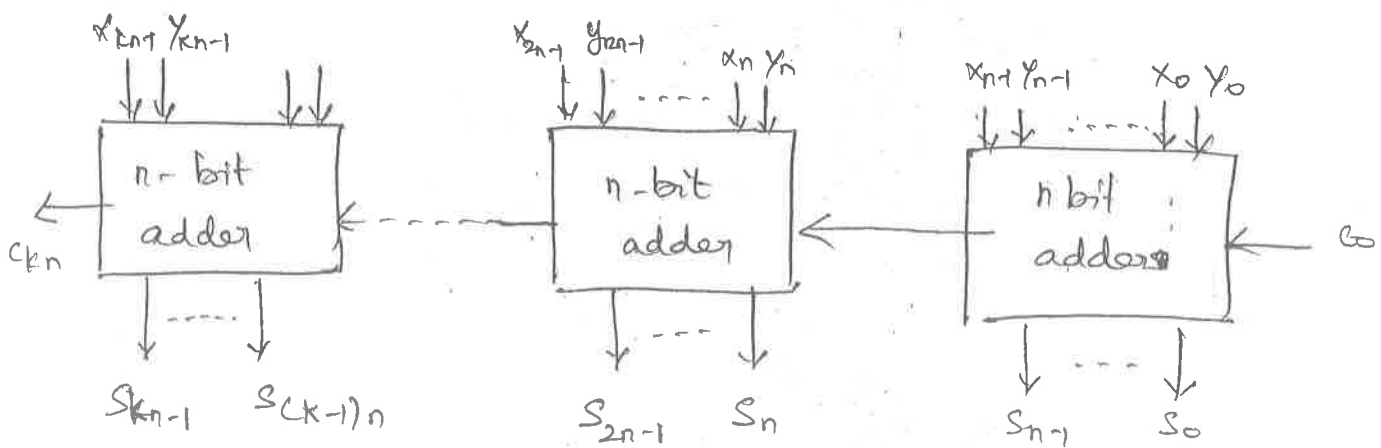
for addition.

The Maximum propagation delay of an n -bit ripple carry adder, the operating speed is nd , where d is the delay of a full adder stage.

Fast Adder Circuit:-

In an n -bit parallel adder (ripple carry) adder there is

The carry signals are also useful for interconnecting k -adders to form an adder capable of handling input numbers that are kn bits long.



Cascade of k n -bit adders.

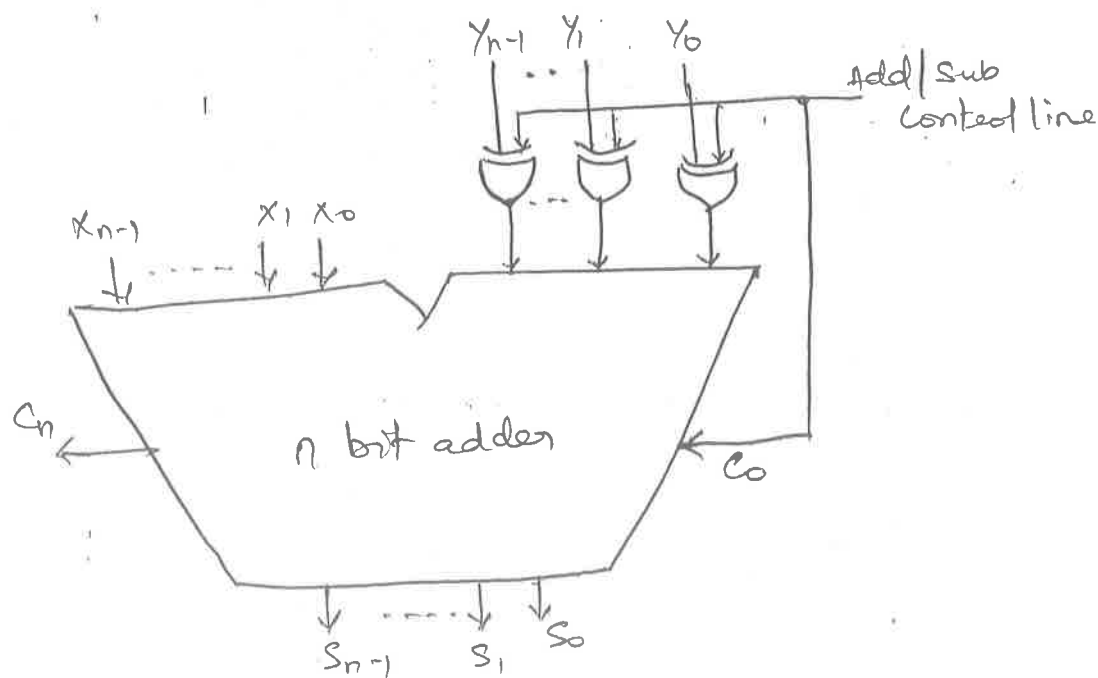
Addition or Subtraction logical Unit:-

The addition / subtraction logic circuit is used to perform either addition or subtraction based on the value applied to Add / Sub input control line.

The line is set to 0 for addition, applying the Y vector unchanged to one of the adder inputs along with a carry-in signal, C_0 , of 0.

When the add/sub control line is set to 1, the Y vector is 1's complemented (that is bit complemented) by XOR gates and C_0 is set to 1 to complete the 2's complementation of Y.

An XOR gate can be added to detect the overflow condition $C_n \oplus C_{n-1}$.



Fast Adder Circuit:-

In an n-bit parallel adder (ripple carry adder), there is too much propagation delay in developing the outputs.

Delay is not acceptable, in comparison with the speed of other processor components and speed of data transfers between registers & cache memories.

The delay through the circuit depends on number of gates in the path from input to outputs. In case of the ripple carry adder, the longest path is from inputs at the LSB position to output at the MSB position.

The propagation delay can be reduced by using a carry look-ahead adder in ALU.

* Carry Look ahead adder.

Carry look ahead Adder :- (Fast Adder) 11

The main purpose of design fast adders is to reduce the time required to 'carry' signals from input to output. One way to compute the input carry needed by stage i directly from a carry like signals obtained from all the preceding stages $i-1, i-2, \dots, 0$, rather than waiting for normal carries to ripple slowly from stage to stage. Adders that use this principle is called Carry look ahead adder. Definition

A n bit carry look ahead adder is formed from n stages, each of which is basically a full adder modified by replacing its carry output line C_i by two auxiliary signal called g_i and p_i called generate & propagate.

A fast adder circuit must speed up the generation of the carry signal. the logical expression for $sum(s_i)$ and carry C_{i+1} (carry out) of stage i is

$$S_i = x_i \oplus y_i \oplus C_i \quad \text{--- (1)}$$

and

$$C_{i+1} = x_i y_i + x_i C_i + y_i C_i \quad \text{--- (2)}$$

factoring (2) eqn



$$C_{i+1} = x_i y_i + C_i (x_i + y_i)$$

We can write

$$C_{i+1} = G_i + P_i C_i \quad \text{--- (3)}$$

Where

$$\begin{aligned} G_i &= x_i y_i \\ P_i &= x_i + y_i \end{aligned}$$

--- (4)

--- (5)

(These two expressions ~~4 & 5~~ is called as Generate and propagate function of stage i.)

If generate function is 1, then $C_{i+1} = 1$ independent of 'C_i' which is the input carry. This occurs when x_i & y_i are 1.

The propagate function will produce output carry when either x_i or y_i is 1.

All G_i & P_i function can be formed independently and in parallel in one logic gate delay after x and y vertices are

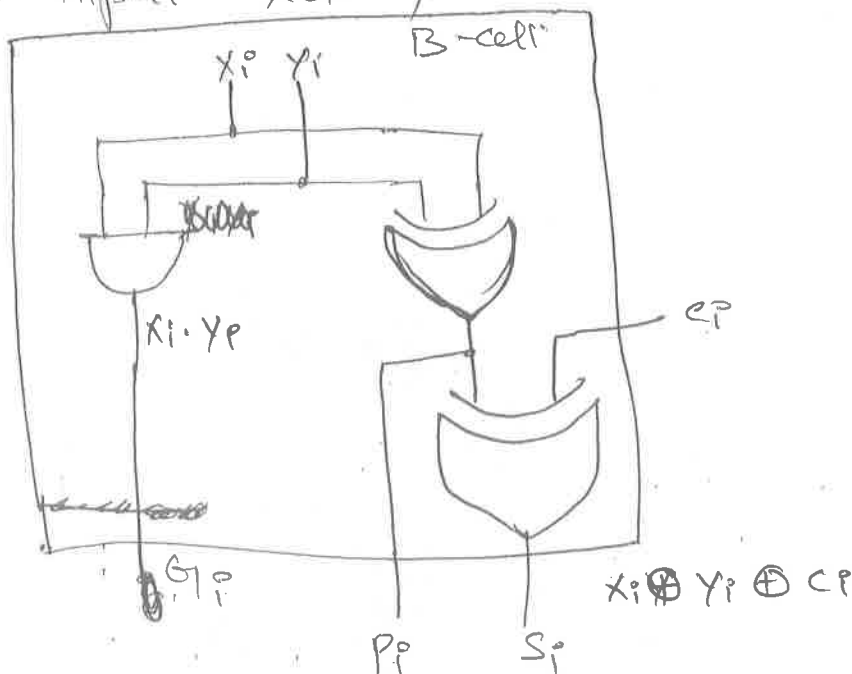
(12)

applied to the inputs of an n-bit adder.

Each Bit Stage contains an AND Gate to form G_i , an OR gate to form P_i and three input XOR gate to form S_i .

A Simple Circuit can be derived called B-cell, where $P_i = x_i \oplus y_i$ which differs \wedge when $x_i = y_i = 1$.

3 input XOR function is realized to two input XOR.



from ② eqn C_P can be obtained, by replacing the subtracting 1 from it

$$C_P = G_{i-1} + P_{i-1} C_{i-1} \quad \text{--- ⑥}$$

Sub (6) in (3) we get

$$C_{i+1} = G_i P + P_i C_i \quad \rightarrow \text{Level 1}$$

$$C_{i+1} = G_i + P_i G_{i-1} + P_i P_{i-1} \underline{C_{i-1}} \quad \text{Level 2} \quad (7)$$

Sub the ~~value~~ ^{expression} of C_{i-1} in (5)

C_{i-1} is obtained from 6 eqn by subtracting 1.

$$\underline{C_{i-1}} = G_{i-2} + P_{i-2} C_{i-2} \quad (8)$$

Sub 8 in (7) we get

$$C_{i+1} = G_i + P_i G_{i-1} + P_i P_{i-1} G_{i-2} + P_i P_{i-1} P_{i-2} \begin{matrix} C_{i-1} \\ C_{i-2} \end{matrix} \quad (9)$$

Continuing expanding the final expression for any carry variable is general form:-

$$C_{i+1} = G_i + P_i G_{i-1} + P_i P_{i-1} G_{i-2} + \dots + P_i P_{i-1} P_{i-2} \dots P_0 C_0 \quad (10)$$

Let Us consider design of 4 bit adder, the carries are implemented as -

Consider eqn (3), (7), (9)

sub $i=0, i=1, i=2, i=3$ respectively

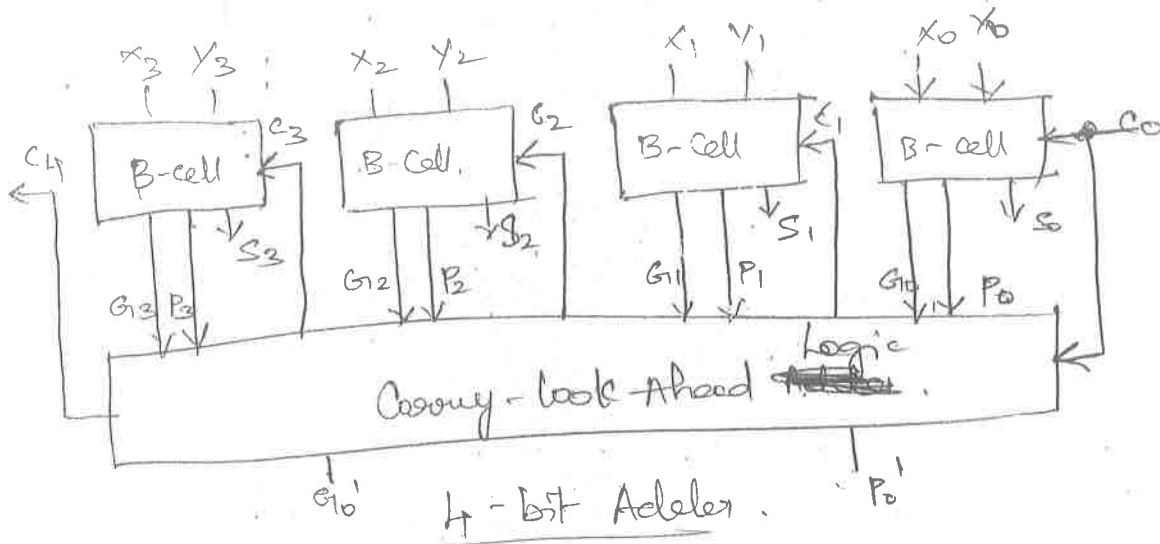
we get

$$i=0 \text{ in } \textcircled{3} \Rightarrow C_1 = G_0 + P_0 C_0$$

$$i=1 \text{ in } \textcircled{7} \Rightarrow C_2 = G_1 + P_1 G_0 + P_1 P_0 C_0$$

$$i=2 \text{ in } \textcircled{9} \Rightarrow C_3 = G_2 + P_2 G_1 + P_2 P_1 G_0 + P_2 P_1 P_0 C_0$$

$$i=3 \quad C_4 = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_0$$



In this diagram the carries are implemented in a block called "Carry Lookahead Logic". An adder implemented in this form is called Carry look ahead adder.

Delay through adder is 3 gate delay for all carry bits and 4 gate delay for all sum bits.

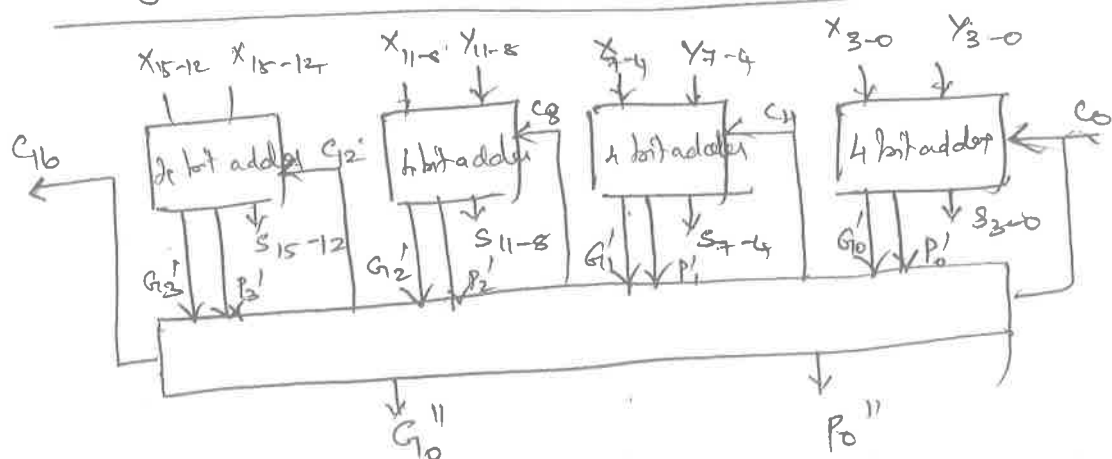
If - we try to extend the carry-lookahead adder for more operands

we run into a problem called gate fan in constraints. From the general expression we see that last AND and OR gate requires 1+2 fan in in generating C_{i+1} . This is a limitation. So the adder cannot be extended directly from the previous B-cells, so it is possible to build longer adders.

~~Warning~~ Fan in :- no of inputs that logic gate can accept. If input exceeds, the output will be undefined.

Eight 4-bit Carry-look ahead adders can be connected to form a 32 bit adder. In order to extend adder, key idea is to generate carries C_4, C_8, \dots in parallel similar way of C_1, C_2, C_3 & C_4 .

Higher Level Generate & propagate functions



(14)

This is a 16 bit adder built from four 4 bit adder blocks, these blocks provide new output function defined as G'_k and P'_k where $k = 0$ for 1st 4 bit block, $k = 1$ for ~~first~~ second 4-bit block and so on.

In the first block,

$$P'_0 = P_3 P_2 P_1 P_0 \quad \text{--- (A)}$$

$$G'_0 = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 \quad \text{--- (B)}$$

\therefore carry C_{16} is formed by one of carry-lookahead circuits as

$$C_{16} = G'_3 + P'_3 G'_2 + P'_3 P'_2 G'_1 + P'_3 P'_2 P'_1 G'_0 + P'_3 P'_2 P'_1 P'_0 G_0$$

\therefore We can extend the operands in the similar way.

A 64 bit adder can be built from four of the 16 bit adders.

\therefore Delay through the carry look ahead adder is 3 gate delay for carry, 4 gate delay for sum, which is less when compared to ripple carry adder.

Multiplication

Multiplication is a slightly more complex operation than addition or subtraction, being implemented by shifting as well as addition. Multiplying two n -bit values together can result in a value of up to $2n$ bits. Because of the partial products involved in most multiplication algorithms more time and more circuit area is required to compute, allocate, and sum the partial product to obtain the multiplication result.

Consider a normal Binary Multiplication

$$\begin{array}{r} \text{Multiplicand} \quad 1000_2 \quad (8_{10}) \\ \text{Multiplier} \quad \quad 1001_2 \quad (9_{10}) \\ \hline 1000 \\ 0000 \\ 0000 \\ 1000 \\ \hline \text{Product} \quad 1001000_2 \quad (72_{10}) \end{array}$$

$$\begin{array}{r} 2 \overline{) 72} \\ 2 \overline{) 36} - 0 \\ 2 \overline{) 18} - 0 \\ 2 \overline{) 9} - 0 \\ 2 \overline{) 4} - 1 \\ 2 \overline{) 2} - 0 \\ 1 - 0 \end{array}$$

The first operand is called the multiplicand, and the second operand is called as the multiplier.

In this example, we restricted ~~as the result~~ of multiplying two 3-bit numbers the decimal digits to 0 and 1. with only two choices, each step of multiplication is simple.

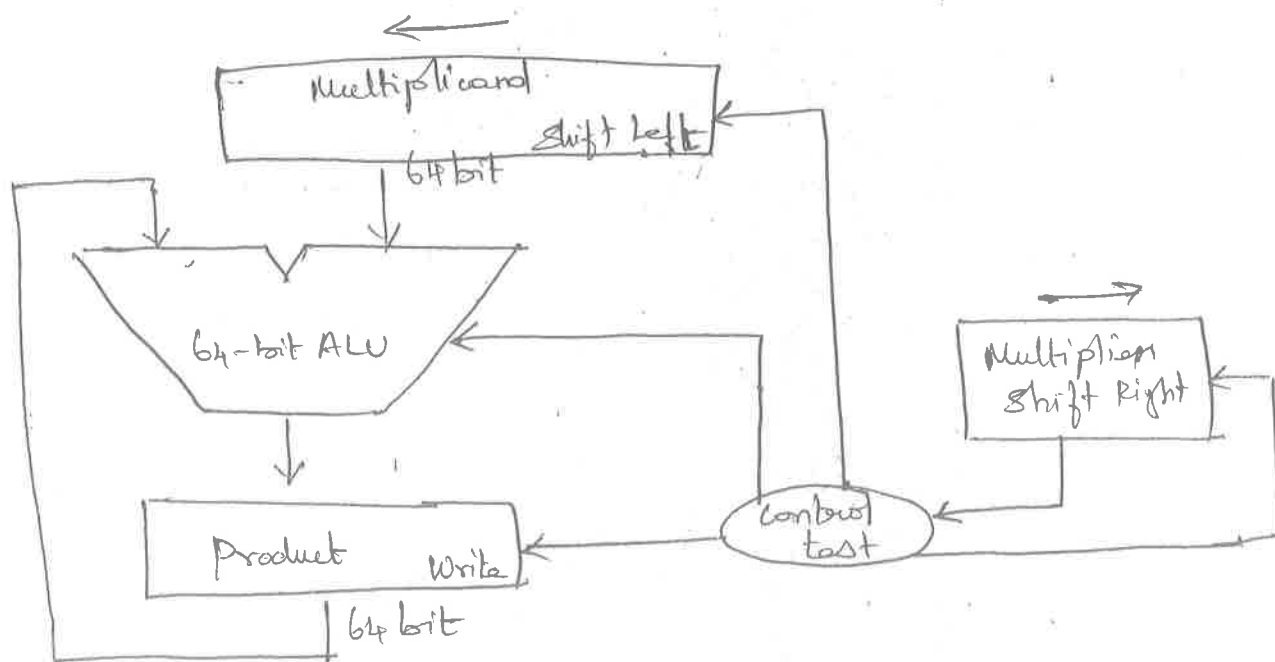
- ① Just place a copy of multiplicand ($1 \times$ multiplicand) in the proper place if the multiplier digit is a 1 (or)
- ② Place 0 ($0 \times$ multiplicand) in the proper place if the digit is 0.

Some of the highly optimized multiplication hardware are:

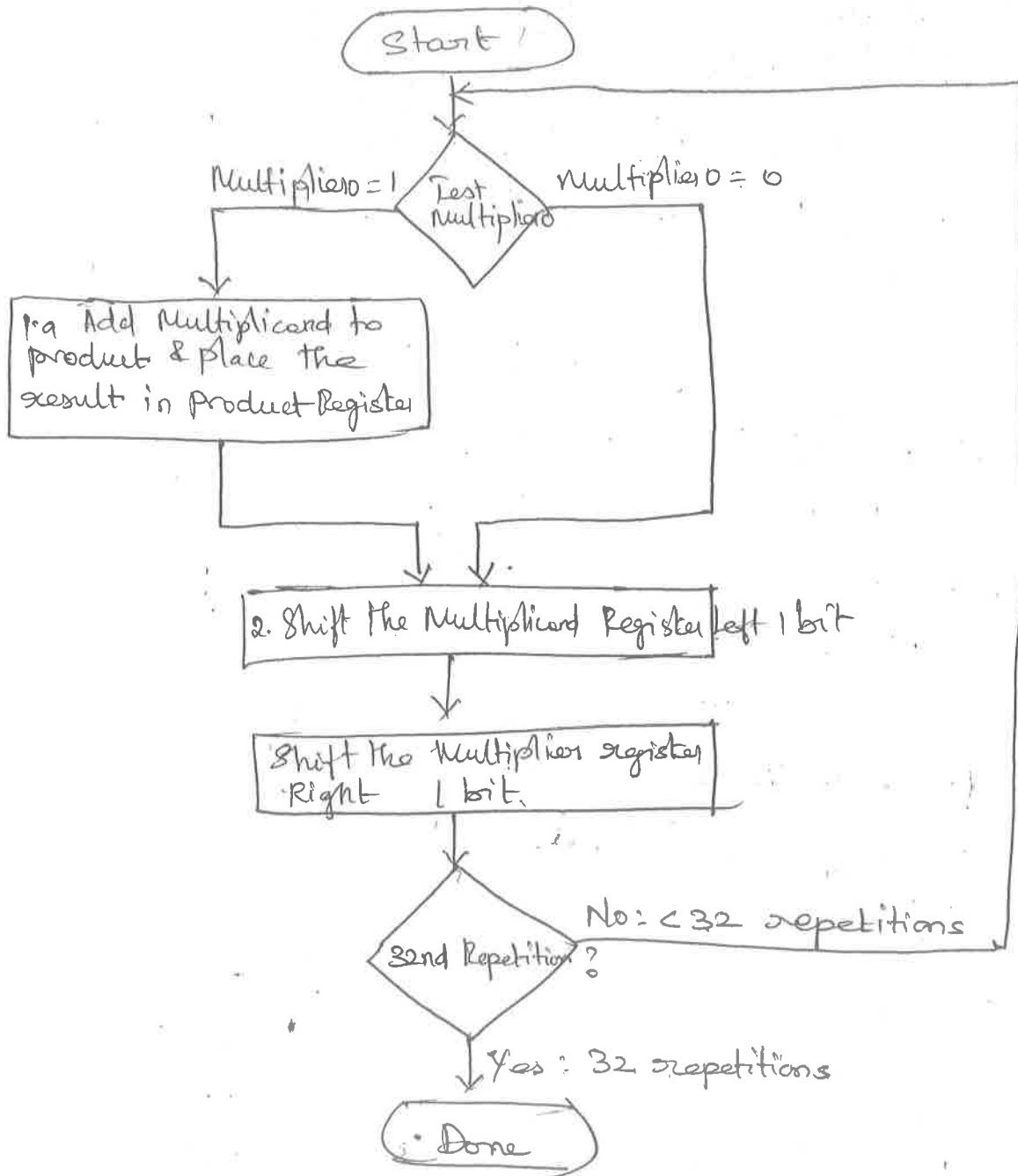
- ① Sequential Version (first version) of the multiplication Algorithm & Hardware
- ② Refined version of the multiplication hardware
- ③ Signed Multiplication \rightarrow (Booth's Algorithm)
- ④ Faster Multiplication \rightarrow carry save adders
- ⑤ Multiply in MIPS.

Sequential Version of the Multiplication Algorithm:-

The hardware is made up of 32 bit multiplier register and 64 bit product register that is initialized to 0. Over 32 steps a 32 bit multiplicand would move 32 bit to the left. Hence we need 64 bit multiplicand register initialized with 32 bit multiplicand in the right half and zero in the left half. The register is then shifted left 1 bit each step to align the multiplicand with the sum being accumulated in the 64 bit product register. The ALU is 64 bit wide to perform 64 bit addition.



The first Multiplication Hardware Algorithm using the hardware.



The three basic steps needed to perform multiplication of each bit (19)

Step 1: The LSB bit of the multiplier determines whether the multiplicand is added to the product register.

Step 2: Shift the multiplicand by 1 bit left, that will have a effect of moving the intermediate operands to the left.

Step 3: Shift the multiplier right by 1 bit that will give the next bit of the multiplier to be examined.

These 3 steps are repeated for 32 times to obtain the product.

eg: Consider 4 bit Multiplication

$$4_{10} \times 3_{10} = 12_{10} \quad \text{or} \quad 0100_2 \times 0011_2$$

Iteration	Steps	Multiplier	Multiplicand	Product
Initiated 0	Initial Values	0010 ^①	0000 0100	0000 0000 +
1.	1. 1 → prod = prod + M _{and} 2. Shift Left Multiplicand 3. Shift Right Multiplier	0011 0011 0001 ^①	0000 0100 0000 1000 0000 1000	0000 0100 0000 010 0000 010
2.	1. 1 → prod = prod + M _{and} 2. Shift Left Multiplicand 3. Shift Right Multiplier	0001 0001 0000 ^①	0000 1000 0001 0000 0001 0000	0000 1100 0000 110 0000 1100
3.	1. 0 ⇒ No operation 2. Shift Left Multiplicand 3. Shift Right Multiplier	0000 0000 0000 ^①	0001 0000 0010 0000 0010 0000	0000 1100 0000 110 0000 1100
4	1. 0 ⇒ No operation 2. Shift Left Multiplicand 3. Shift Right Multiplier	0000 0000 0000	0001 0000 0100 0000 0100 0000	0000 110 0000 110 0000 110

It is observed that

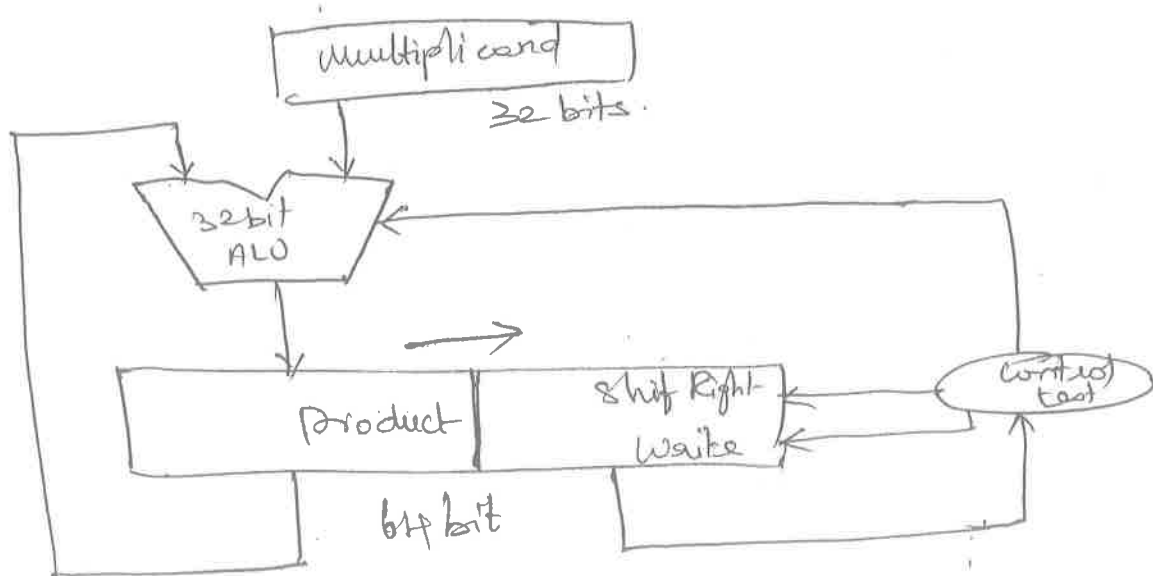
- ① Each step took a clock cycle.
- ② Half bits in Multiplicand is always 0.
∴ 64-bit adder is wasted.
- ③ LSB of the product is never changed once formed, instead of shifting multiplicand to left, shift product to right:-

(18)

These drawbacks were considered for the another version of the multiplier

Refined Version of the Multiplication Hardware:-

In comparison to the first Version, the Multiplicand Register, ALU and Multiplier Register are all 32 bit wide with only the product register is 64 bit. The product is shifted right. The separate Multiplier is disappeared and it is placed in the Right half of the product register.



Faster Multiplication! - → Fast Multiplier Hardware
→ Carry same address

→ Moore's law provided so much more in resources that hardware designers can now build much faster multiplication hardware.

→ Rather than to use a single 32-bit address 31 times, the hardware "unrolls the loop" is it uses 32 bit address for each bit of multiplier and organizes them in order to minimize the delay.

→ It provides a 32 bit address for each bit of multiplier, one input is multiplied ANDed with a multiplier bit and other is the output of a prior address.

→ Straight forward approach is to connect the outputs of the address on the right, to the inputs of the address on the left. making a stack of address 32 bit high.

→ An alternate way to organize these 32 additions in a "parallel tree" as shown in

Faster Multiplication:- → Fast Multiplier Hardware
→ Carry same adders

→ Moore's law provided so much more in resources that hardware designers can now build much faster multiplication hardware.

→ Rather than to use a single 32-bit adder 31 times, the hardware "unrolls the loop" is it uses 32 bit adder for each bit of multiplier and organizes them in order to minimize the delay.

→ It provides a 32 bit adder for each bit of multiplier, one input is multiplied ANDed with a multiplier bit and other is the output of a prior adder.

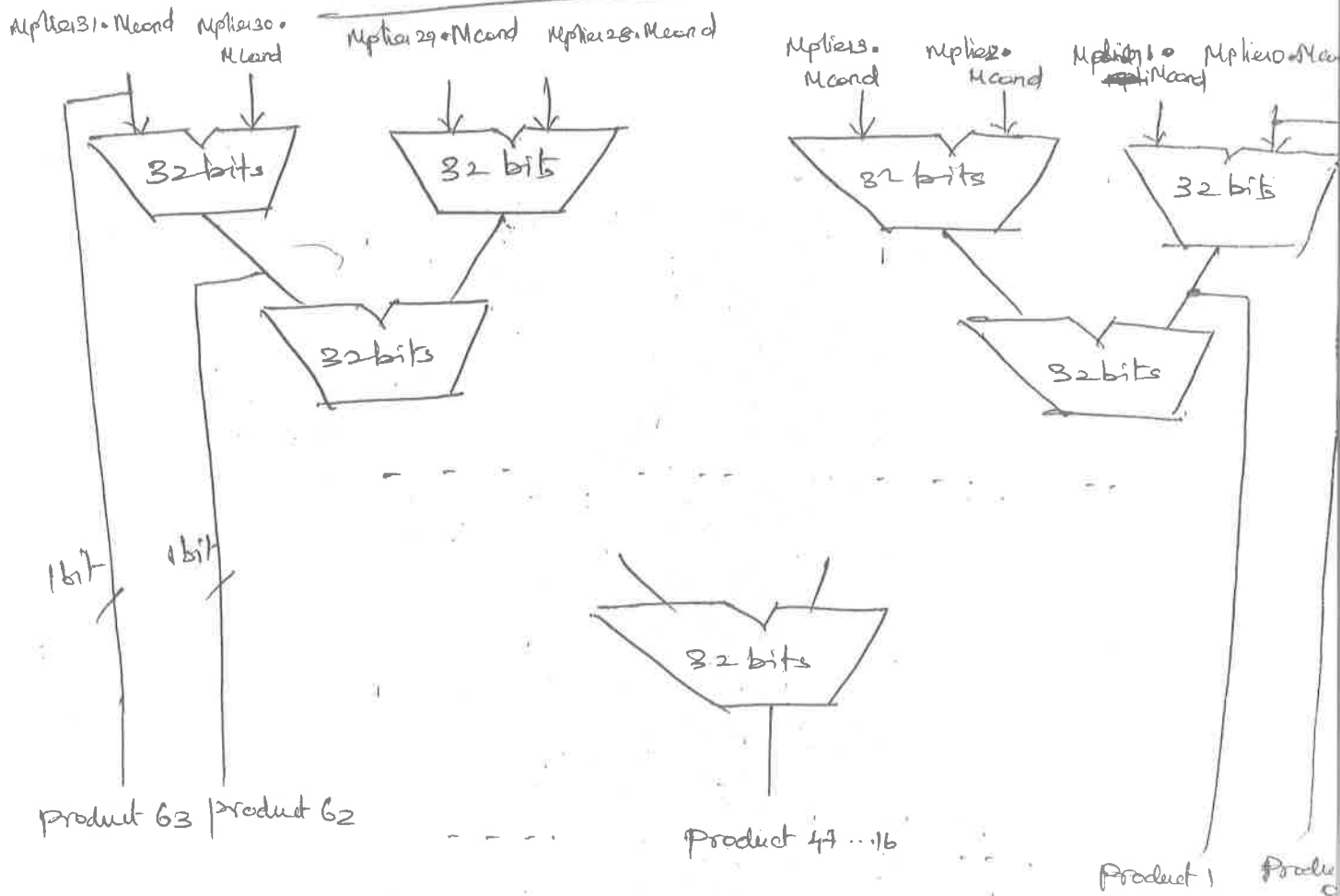
→ Straight forward approach is to connect the outputs of the adders on the right, to the inputs of the adders on the left, making a stack of adders 32 bit high.

→ An alternate way to organize these 32 additions in a "parallel tree" as shown in

the diagram.

19

Fast Multiplication Hardware



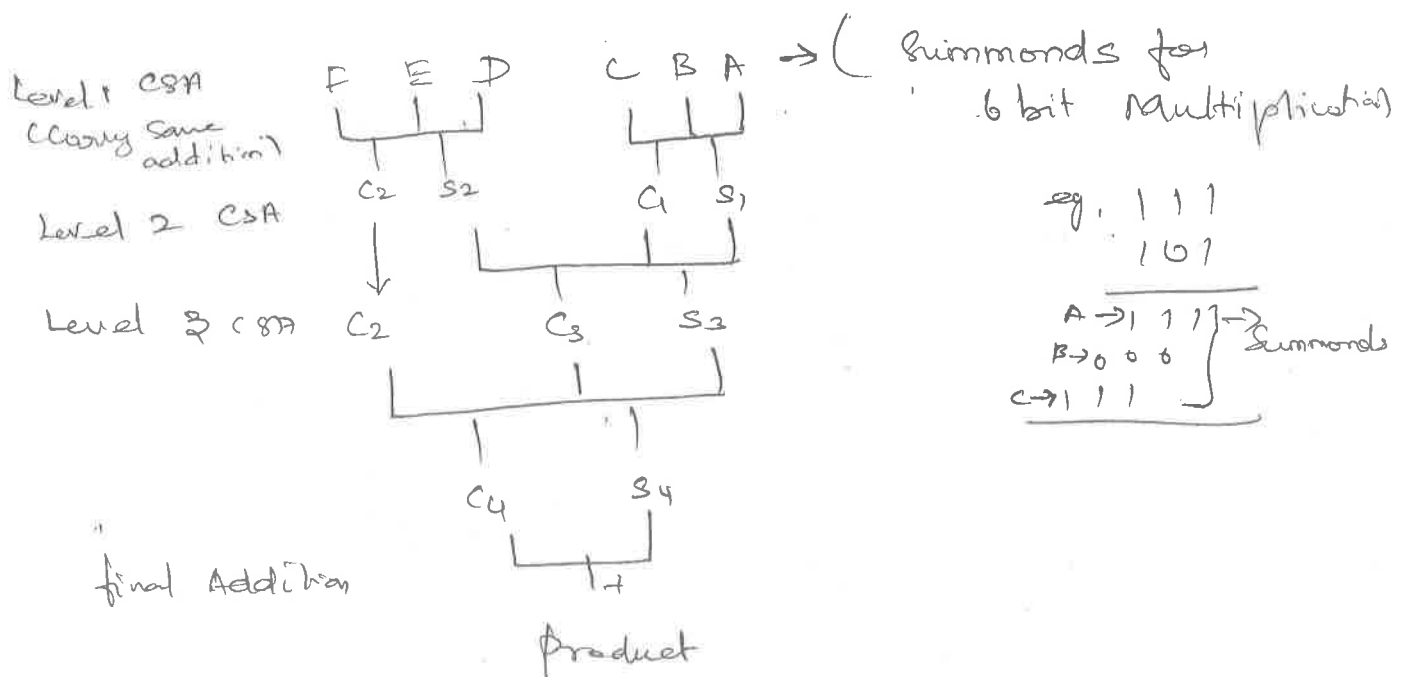
→ Instead of waiting 32 add times, we just wait $\log_2(32)$ or 5 fine 32-bit add times.

* Multiplication can go even faster than fine add times because of the use of "carry save adders"

Carry Save Adders:- (for Multiplication)

"A technique called carry save addition (CSA) speeds up the addition process by letting the saving the carry and introducing it in the next row instead of letting the carry to ripple along the "row".

Schematic Representation of The carry save addition



Signed Multiplication: (Booths Algorithm)

In signed number there is a need to take the sign of the number into consideration. The Easiest way to deal with signed number is first convert the multiplier and multiplicand to positive and then remember the original signs. The algorithm should then run for 32 iterations leaving the signs out of calculation.

→ We need to negate the product only if the original signs disagree (ie signs differ).

eg: consider the case of a [ⓐ]positive multiplier & negative multiplicand, when we add a negative multiplicand to a partial product we must extend the sign-bit value of each multiplicand to the left as far as the product will extend.

eg: $-4 \times 3 \Rightarrow 1100_{10} \times 0011_{10}$

1100 (-4)
x 0011 (3)

Sign bit extension ←

```

  1111 1100
  1111 1000
  0000 00
  0000 0
  -----
  1111 0100 ⇒ -1210
  
```

```

  4 → 0100
  15 → 1011
  +1 → 1
  -----
  1100
  12 → 01100
  15 → 10011
  " 1
  -----
  10100
  
```


For negative ^② Multiplier, a straight forward solution is to form the 2's complement of both the multiplier and the multiplicand and proceed as in the case of a positive multiplier.

BOOTH'S ALGORITHM:-

In order to take care of both positive and negative multipliers Booth's Algorithm can be used.

The Booth's Algorithm generates a 2n-bit product and treats both positive and negative ~~also~~ 2's complement uniformly.

→ In the Booth's Scheme, booth's recoding of a multiplier should be known.

Multiplier		Variation of Multiplicand Selected by bit i
Bit i	Bit i-1	
0	0	$0 \times M$
0	1	$+1 \times M$
1	0	$-1 \times M$
1	1	$0 \times M$

eg: Booth's recoding for the multiplier.

0010110100110110 [0] ^{called implied 0}

0+1 -1 +1 0 -1 +1 -1 0 +1 0 -1 +1 0 -1 0

↳ Always assume i-1 bit in initial stage is 0

↓

Called as the Booths recoding.

⇒ Note: Join the two bits from Right to Left then look the Recoding table for the two bit combination if 01 [0] the 10 in table will correspond to -1XM is the recoding for 10.

Different:-

Cases of Multiplier :-

- 1) Worst Case Multiplier

01 01 01 01 01 01 01 01 [0]

+1 -1 +1 -1 +1 -1 +1 -1 +1 -1 +1 -1 +1 -1
- 2) Ordinary Case Multiplier

1 1 0 0 0 1 0 1 1 0 1 1 1 (0)

0 -1 0 0 +1 -1 +1 0 -1 +1 0 0 -1
- 3) Good Multiplier

0 0 0 0 1 1 1 1 0 0 0 (0)

0 0 0 +1 0 0 0 0 -1 0 0 0

eg: Recode the Multiplier 101100 for Booth's Multiplication

Multiplier 1 0 1 1 0 0 [0] → Implied Zero
 Recoded Multiplier $\boxed{-1 +1 0 -1 0 0}$

"The speeding up of Multiplication is achieved depending upon the multiplier. If the multiplier has few contiguous blocks of 1's then summands will be reduced, the speedup will be more."

eg: Note ⇒ Always record the multiplier (if it +ve or Negative)

1) If the Multiplier is -ve and Multiplicand is +ve.

$$6_{10} \times -5_{10} = -30$$

$$6_{10} \rightarrow 0110$$

$$-5 \rightarrow 1011$$

$$\begin{array}{r} -5 \rightarrow 0101 \\ 1010 \\ \hline 1011 \end{array}$$

Solution:-

$\begin{pmatrix} 0 & 1 & 1 & 0 & (+6_{10}) & \text{(Multiplicand)} \\ 1 & 0 & 1 & 1 & (-5_{10}) & \text{(Multiplier)} \\ -1 & +1 & 0 & -1 & & \text{(recoded Multiplier)} \end{pmatrix}$

Take the multiplicand and regarded multiplies and perform multiplication.

$$\begin{array}{r}
 0 1 1 0 \\
 +1 +1 0 -1 \\
 \hline
 1 1 (0) 1 1 1 0 1 0 \\
 0 0 0 0 0 0 0 \\
 0 0 0 1 1 0 \\
 \oplus 1 1 0 1 0
 \end{array}$$

$$\begin{array}{r}
 0110 \\
 1001 \\
 \hline
 1010
 \end{array}$$

$$\begin{array}{r}
 * 1 1 1 0 0 0 1 0 \Rightarrow (-30) \\
 \hline
 \text{discarded}
 \end{array}$$

check:- $30 \rightarrow 00011110$
 1's comp 11100001
 $\underline{11100010} \rightarrow (-30)$

ii) If the multiplicand is +ve and multiplier is positive.

$$\begin{array}{cc}
 -6 & \times & 5 & = & -30 \\
 10 & & 10 & &
 \end{array}$$

$$\begin{array}{l}
 -6 \rightarrow 1010_2 \\
 5 \rightarrow 0101_2
 \end{array}$$

$$\begin{array}{r}
 0110 \\
 1001 \\
 \hline
 1010
 \end{array}$$

Recorded Multiplier (1010) (-6) (Multiplicand)
 (0101) (5) (Recorded Multiplier)
 +1 -1 +1 -1

1010 (-6)
 +1 -1 +1 -1
 0000 0110
 1111 010
 0001 10
 11010
 *11100010₂ (-30)
 → discard

ii) If both the multiplier and multiplicand are negative:

$$-6_{10} \times -5_{10} = +30_{10}$$

-6 → 1010 (Multiplicand)

-5 → (0101) (Multiplier)

→ -1 +1 0 -1

Booths
 recording of
 Multiplier

$$\begin{array}{r} -6) 0110 \\ \underline{1001} \\ 51010 \\ -5) \\ \underline{0101} \\ 1010 \\ \underline{1011} \end{array}$$

1010
-1+10-1

(-6)₁₀

(Booths Recoding of Multiplier)

0000 0110

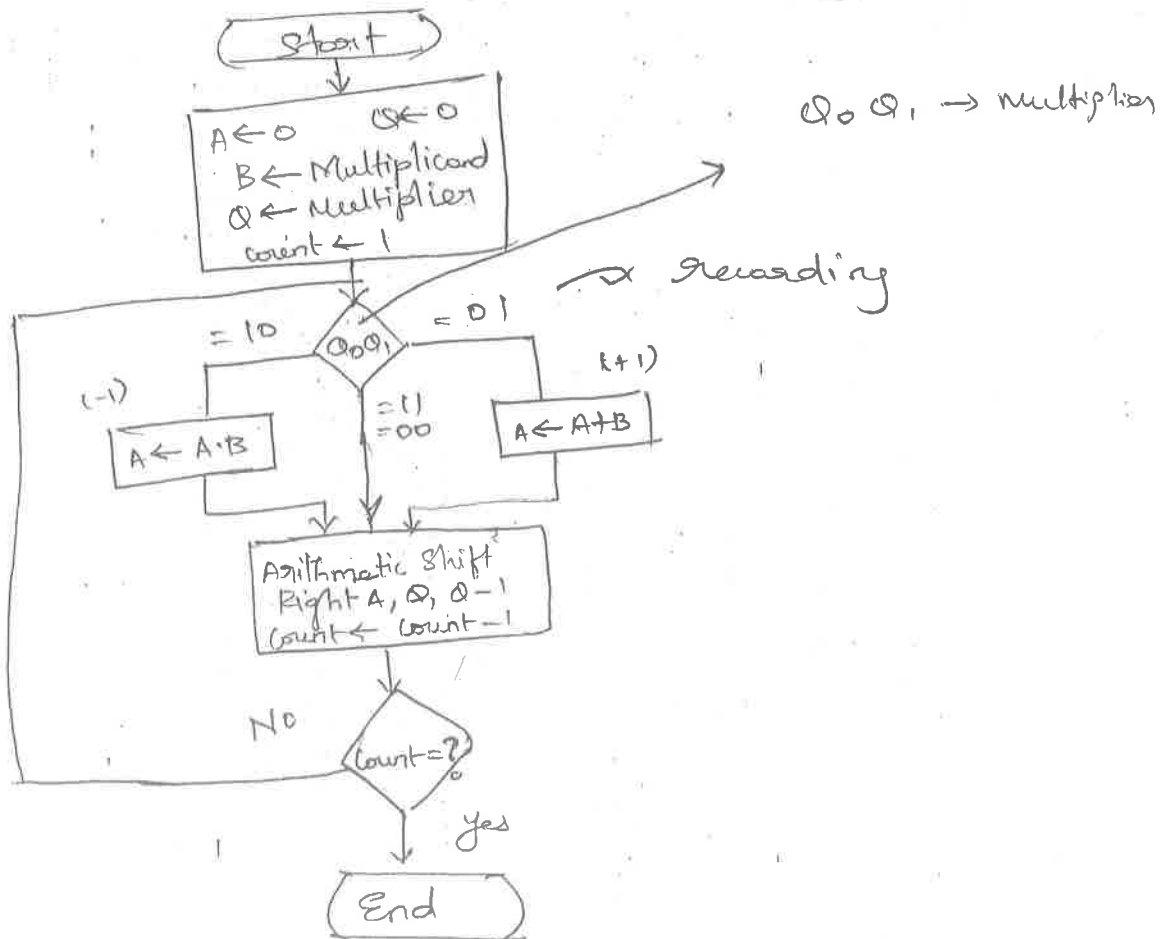
0000 000

1110 10

(+) 00110

~~1~~ 0001110₂ = (+30)₁₀
↓ discard

Flow Chart Booths Algorithm for Signed Multiplication



Features of Booths Algorithm:-

* It handles both positive and negative multipliers uniformly.

* Second It achieves some efficiency in the number of additions required when the multiplier has a few large blocks of 1's.

Multiply in MIPS:-

MIPS has two instructions mult (multiply) and Multiply Unsigned (multu). MIPS provide a separate pair of 32 bit registers to hold 64 bit product.

Division

- ① Division
- ② First Version of the Division Hardware & Algorithm
- ③ Refined Version of the Division Hardware
- ④ Signed Division
- ⑤ Faster Division
- ⑥ Division in MIPS

Division:-

→ The reciprocal operation of multiplication is divide.

→ It is very complex similar to multiplication

→ Division can be achieved by successive subtraction.

eg:- Consider the binary division of 74_{10} by 8_{10}

binary form is 1001010_2 by 1000_2

$$\begin{array}{r}
 1001010 \\
 \underline{1000} \\
 1010 \\
 \underline{1000} \\
 0010
 \end{array}$$

1000 \swarrow Divisor
 1001010 \rightarrow Dividend
 1010
 0010 \rightarrow Remainder

There are two operands called, the dividend and the divisor. The result is called as the Quotient, second result is called remainder.

dividend - A number being divided.

divisor - A number that the dividend is divided by.

Quotient - The primary result of a division.

remainder - The secondary result of a division.

We can express the relation between these by,

$$\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}$$

where the remainder should be smaller than that of the divisor.

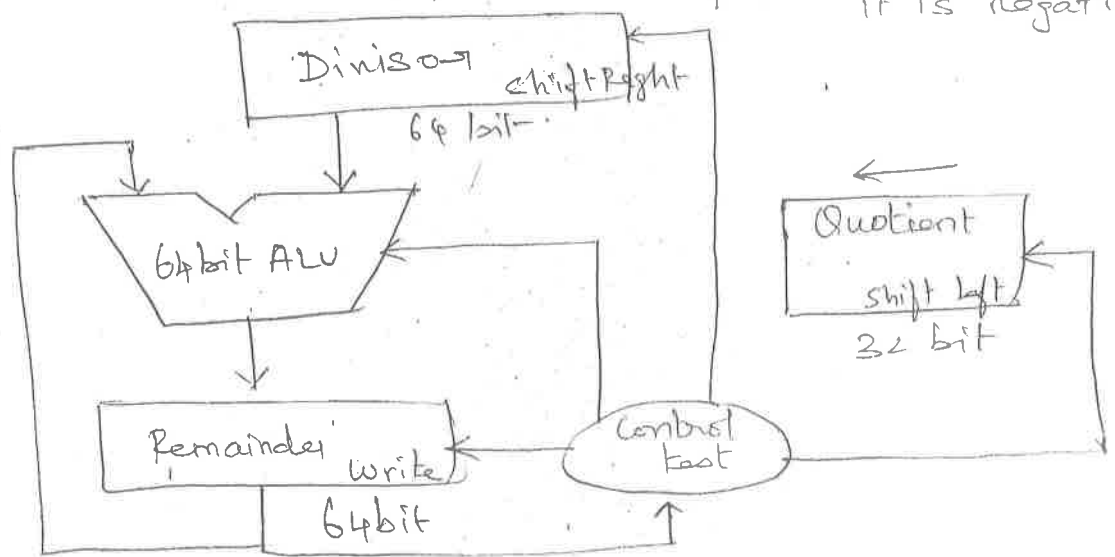
First Version of the Division Algorithm and Hardware :

Let us assume that the dividend and the divisor is positive, Hence the Quotient and the remainder is non-negative. The Hardware uses

32 bit Quotient Register and set to 0. It uses

64 bit ALU.

Each iteration of the algorithm needs to move the divisor to the right on digit, so the divisor is placed in the left half of 64 bit divisor register and shift it right 1 bit each step. The Remainder Register is initialised with the dividend. It is also called Restoring division because the remainder is restored when it is negative.



First Version of the division hardware

Steps involved in division Algorithm:-

The computer isn't smart enough to know in advance whether the divisor is smaller than the dividend.

Step 1: Test $\text{divisor} < \text{dividend}$ (ie) Subtract the divisor from the dividend and place the result in the remainder register.

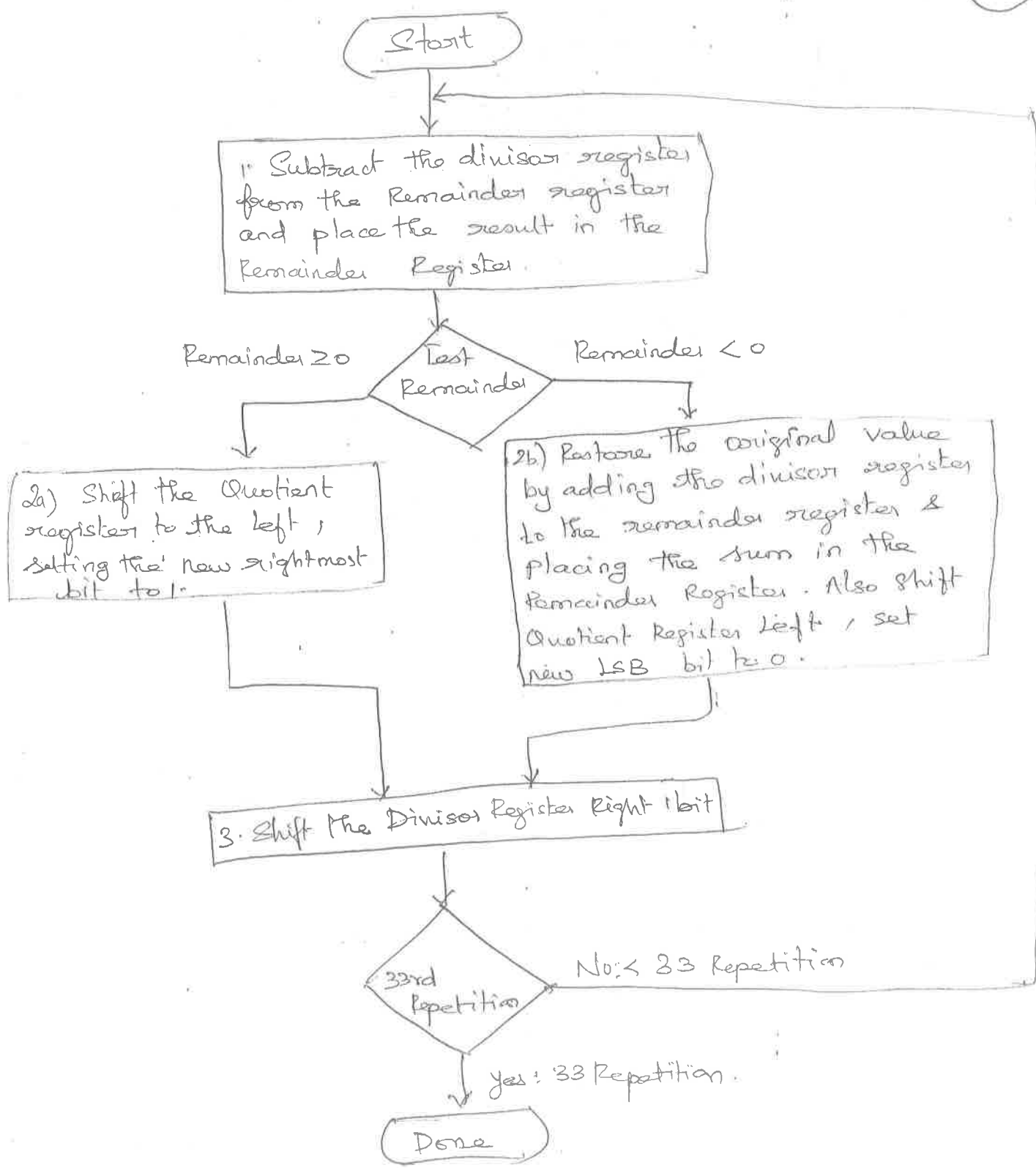
Step 2: a) if $\text{remainder} \geq 0$ $\text{divisor} < \text{dividend}$; shift the quotient register to the left setting the new right most bit to 1.

b) if $\text{remainder} < 0$, Restore the original value by adding the divisor to the dividend (Remainder Register) and placing the sum result in remainder register and also shift the quotient register to left, setting the new least significant bit to '0'.

Step 3: shift the divisor right by 1 bit.

These 3 steps are repeated for < 32 times to obtain the remainder and Quotient result.

The division algorithm, using the hardware is given below.



A division Algorithm, Using The hardware.

eg: divide 7_{10} by 2_{10}
 the binary is $0000\ 0111_2$ by 0010_2

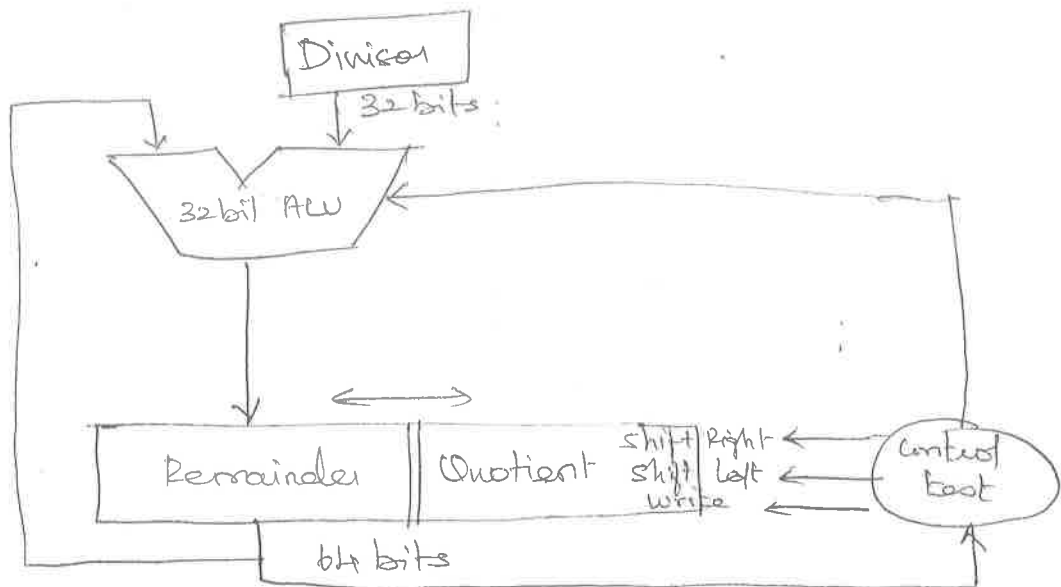
Iteration	Steps	Quotient	Divisor	Remainder
0	Initial Values	0000	0010 0000	0000 0111
1	1. Rem = Rem - Div 2. Rem < 0 \Rightarrow +Div, sll Q, Q ₀ =0 3. Shift Div Right	0000 0000 0000	0010 0000 0010 0000 0001 0000	0110 0111 0000 0111 0000 0111
2	1. Rem = Rem - Div 2. Rem < 0 \Rightarrow +Div, sll Q, Q ₀ =0 3. Shift Div Right	0000 0000 0000	0001 0000 0001 0000 0000 1000	0111 0111 0000 0111 0000 0111
3	1. Rem = Rem - Div 2. Rem < 0 \Rightarrow +Div, sll Q, Q ₀ =0 3. Shift Div Right	0000 0000 0000	0000 1000 0000 1000 0000 0100	0111 1111 0000 0111 0000 0111
4	1. Rem = Rem - Div 2. Rem \geq 0 \Rightarrow sll Q, Q ₀ =1 3. Shift Div Right	0000 0001 0001	0000 0100 0000 0100 0000 0010	0000 0011 0000 0011 0000 0011
5	1. Rem = Rem - Div 2. Rem \geq 0 \Rightarrow sll Q, Q ₀ =1 3. Shift Div Right	0000 0011 0011	0000 0010 0000 0010 0000 0001	0000 0001 0000 0001 0000 0001

A revised / Improved Version of Division Hardware :-

\rightarrow The Divisor Register, ALU and Quotient Register are all 32 bits wide

\rightarrow The Remainder Register alone is left as 64 bits

\rightarrow This Version also combines the Quotient register with the right half of the Remainder Register.



An Improved Version of the Division Hardware.

Signed Division:-

Remember the sign of the divisor and dividend and then negate the Quotient if signs disagree (ie if they are not same).

$$\text{eg: } -6_{10} \div 3_{10} = -2_{10} \quad (\text{since diff sign})$$

The one complication is that sign of the remainder.

$$\boxed{\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}}$$

eg: Consider the combinations of $\pm 7_{10}$ by $\pm 2_{10}$ to find the sign.

(i) $+7 \div +2$: Quotient = +3 , Remainder = +1

checking:-

$$7 = 3 \times 2 + (+1) = 6 + 1.$$

(ii) $+7 \div +2$: Quotient = -3,

Remainder can be calculated Using this

$$\begin{aligned}\text{Rem} &= \text{Divident} - (\text{Quotient} \times \text{Divisor}) \\ &= -7 - (-3 \times +2) = -7 - (-6) = -1\end{aligned}$$

So

$$-7 \div +2 \text{ Quotient} = -3, \text{ Remainder} = -1$$

here, -4 and Remainder +1 also fits this formula.

this anomalous behaviour is avoided by following the rule that the divident & Remainder must have same signs. no matter what the signs of the divisor or Quotient is.

$$\underline{+7 \div -2} : \text{Quotient} = -3 \quad \underline{\text{Remainder} = +1}$$

$$\underline{-7 \div -2} : \text{Quotient} = +3, \quad \underline{\text{Remainder} = -1}$$

Faster Division:-

-> We can use adders to speed up multiply, but we cannot use the same for division.

The reason is that we need to know the signs of the difference before performing the next step of the Algorithm.

→ There are techniques to produce more than one bit of the quotient per step.

(Sweeney, Robertson & Tocher)
↑
creator names

→ The SRT technique tries to guess several 4bit quotient bits per step using a table lookup based on the upper bits of the dividend and remainder.

→ The accuracy of this fast method depends on having proper values in the lookup table.

→ This technique tries to guess several 4bit quotient bits per step, using the lookup table based on upper bits of dividend and remainder.

→ These algorithms use 6 bits from the remainder ~~here~~ and 4 bits from divisor to index a table.

Division in MIPS:

* MIPS supports multiplication and division using existing hardware, ~~the~~ ALU & shifters.

→ It needs one extra hardware component - a 64 bit register able to support sll and ~~sra~~ sra instructions.

→ Signed division is supported by 'div' instruction and unsigned is by 'divu' instruction.

→ MIPS hardware doesn't check for division by zero. Thus 'divide by zero' exception must be handled in system software.

Non Restoring Division:-

⊗^{2m} (The restoring division algorithm can be improved by avoiding the need for restoring the Remainder Register 'A' after an unsuccessful subtraction.)

(Subtraction is unsuccessful if the result is negative.)

Remainder is in 'A' register.

nbit, Divisor is loaded in M register.

nbit dividend is loaded into Register Q at the start of the operation.

Non Restoring division Algorithm:-

① Step 1:- Do the following n times:

1. If the sign of A is 0, shift A and Q left one bit position and subtract M from A otherwise shift A and Q left and add M to A.
2. Now if the sign of A is 0, set q_0 to 1 otherwise set q_0 to 0.

Step 2:- If the sign of A is 1, Add M to A.

∴ Step 2 is needed to leave the proper positive remainder in A at the end of the n cycles of step 1.

eg: $8_{10} \div 3_{10}$
 $1000_2 \div 0011_2$

Non-Restoring division

Initial Step	<div>A</div> <div>00000 00011 → M</div>	Q	1000
Shift A & Q left	<div>00001 A</div> <div>(-) 00011 M</div> <div>①1110</div>	000□	Step 1
Sub M from A Set q ₀		0000'	(Step 2)
Shift A & Q left	<div>11100</div> <div>(+) 00011</div> <div>①1111</div>	000□	Step 1
Add M to A Set q ₀		0000	(Step 2)
Shift A & Q left	<div>11110</div> <div>(+) 00011</div> <div>①0001</div>	000□	Step 1
Add M		0001	Step 2
Set q ₀			
Shift A & Q left	<div>00010</div> <div>(-) 00011</div> <div>11111</div>	001□	Step 1
Sub M		0010	Step 2
Set q ₀			
Shift A & Q left	<div>00010</div> <div>(-) 00011</div> <div>11111</div>	001□	Step 1
Sub M		0010	Step 2
Set q ₀			

Finally Restore Remainder

$$\begin{array}{r} 1111 \\ 00011 \\ \hline 00010 \rightarrow \text{Remainder} \\ \hline \end{array}$$

Floating Point

29

Floating point representations of decimal numbers are essential to scientific computation using scientific notation with a fraction (F)

an exponent (E) and certain radix (r) in the form $F \times r^E$

Decimal numbers use radix of 10

eg: 1.2×10^{-2}

Binary numbers use radix of 2

eg: 1.1×2^1

Representation of floating point is not unique eg: 55.66 can be represented as

1. 5.566×10^1

2. 0.5566×10^2

3. 0.05566×10^3 and so on.

Some examples of ~~floating point~~ ~~fractions are~~ decimal point numbers are

$3.14159 \dots_{10}$ (π)

0.000000001_{10}

$3,155760000_{10}$

$\Rightarrow \begin{cases} 1.0 \times 10^{-9} \\ 3.15576 \times 10^9 \end{cases}$

Scientific Notation

Scientific Notation - It has a single digit to the left of the decimal point.

$$\boxed{\text{eg: } 3.1459_{10} \times 10^3}$$

Normalized Notation:- A number is scientific and has no leading 0's is called Normalized number.

$$\text{eg: } \boxed{1.0_{10} \times 10^{-9} \rightarrow (\text{Normalized})}$$

$$10.0 \times 10^{-10} \text{ \& } 0.1 \times 10^{-8}$$

(Not Normalized)

Binary numbers can also be represented in Scientific notation

$$\text{eg: } 1.01010_2 \times \underbrace{(2)^3}_{\text{base 2}} \rightarrow \text{base 2.}$$

Floating point:-

It represents numbers in which the binary point is not fixed (because they are shifted to have one non-zero digit to left) they are floating hence it is called floating point.

In binary the form is

$$\boxed{1.\text{XXXXX}_2 \times 2^{\text{YYYY}}}$$

2:00

20

Scientific notation for goals is normalized form offers 3 advantages:-

- ① It simplifies exchange of data that include floating + point numbers.
- ② It simplifies floating point arithmetic algorithms to know that numbers will always be in this form.
- ③ It increases the accuracy of the numbers that can be stored in a word.

Floating-point Representation:-

A designer of a floating point representation must find compromise between the size of the fraction and the size of the exponent

There are three fields in floating-point representation

- ① Fraction - Value generally between 0 and 1
- ② Exponent - It is numerical representation of floating point
- ③ Sign - This indicates the sign of the number if it is positive the bit is '0' if negative it is set to 1.

(ie) In 32-bit Representation:-
 Here S the sign of floating point is 1 bit.
 Exponent - Value including the sign of
 the exponent represented in
 8 bit.
 Fraction :- 23 bit number.

The representation used for Floating
 point is called "Sign & Magnitude"
 representation since sign is a separate
 bit.

There are two representation:-

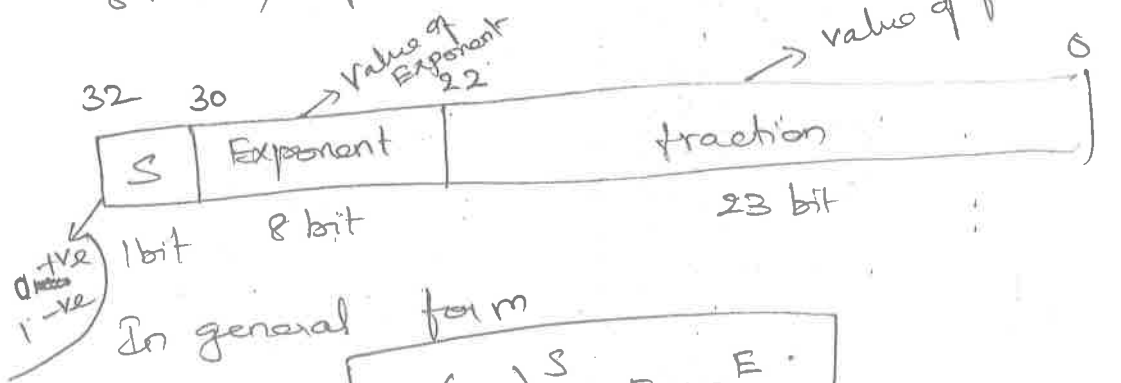
- ① Single precision
- ② Double Precision.

①

Single Precision:-

It uses a 32 bit word to
 represent a floating point Value.

Where S uses 1 bit, exponent uses
 8 bit, fraction uses 23 bit.



$$(-1)^S \times F \times 2^E$$

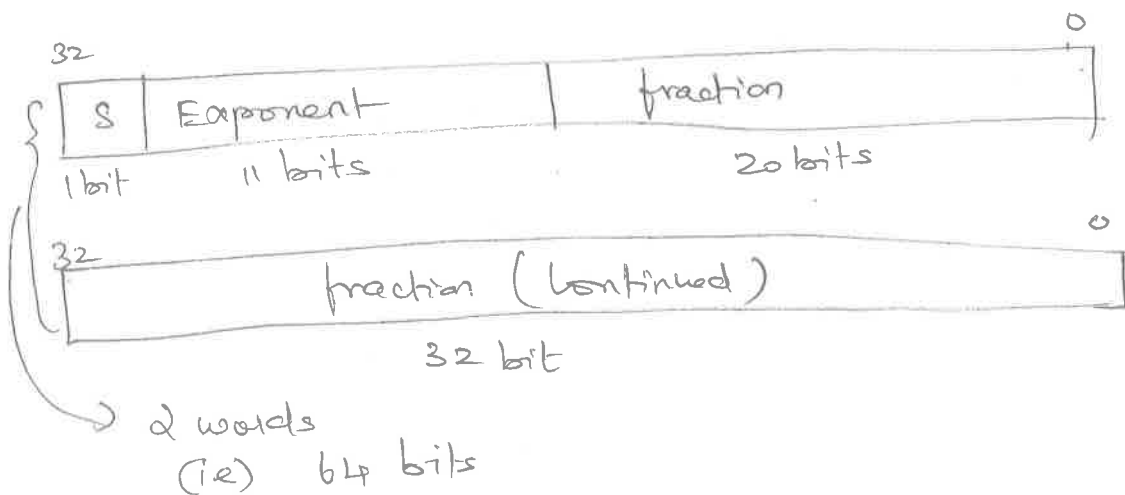
Fractions almost as small as $2.0_{10} \times 10^{-38}$ and numbers almost as large as $2.0_{10} \times 10^{38}$ can be represented in computer.

⇒ Overflow & underflow Interrupt can occur in the floating point arithmetic

Overflow: A situation which positive exponent becomes too large to fit in the exponent field.

Underflow: A situation which negative exponent becomes too large to fit in the exponent field.

One way to reduce Overflow & underflow is to offer format that has larger exponent.
So double precision is used eg: like
double ~~the~~ datatype in C



The double floating point number takes two words. S is sign bit, exponent value of 11 bit exponent field, fraction is the 52 bit number in the fractional field.

It allows numbers almost as small as $2.0_{10} \times 10^{-308}$ and almost as large as $2.0_{10} \times 10^{308}$.

IEEE 754 Standard:-

To pack even more bits into the significant, it makes the leading 1-bit of normalized binary number implicit.

Hence number is 24 bit long in single precision (implied 1 & 23 bit fraction)

Similarly for double precision is 53 bit (1 implied, 52 fraction bits).

Thus $00 \dots 00_2$ represent 0, the representation except for 0 is

$$(-1)^S \times (1 + \text{fraction}) \times 2^E$$