

3. Compare Arithmetic

I. Data Representation.

) Basic Formats

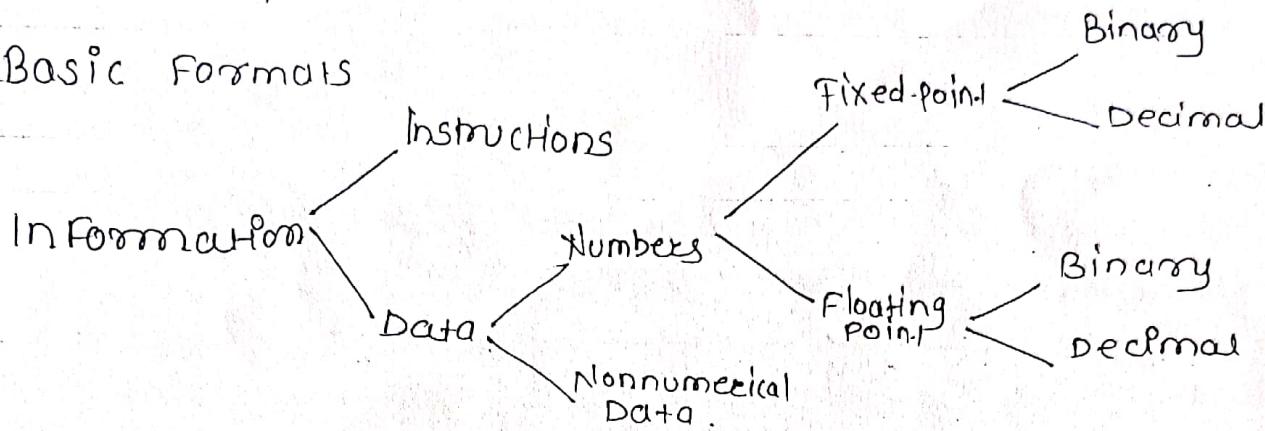


Fig. The Basic information types.

→ Word length.

- Information is represented in a digital computer by means of binary words, where word is a unit of information of some fixed length n . An n -bit word allows up to 2^n different items to be represented.

Example, with $n=4$, we can encode the 10 decimal digits as follows:

$0 = 0000$	$1 = 0001$	$2 = 0010$	$3 = 0011$	$4 = 0100$
$5 = 0101$	$6 = 0110$	$7 = 0111$	$8 = 1000$	$9 = 1001$

- To encode alphanumeric symbols / characters , 8-bit words called bytes are commonly used .

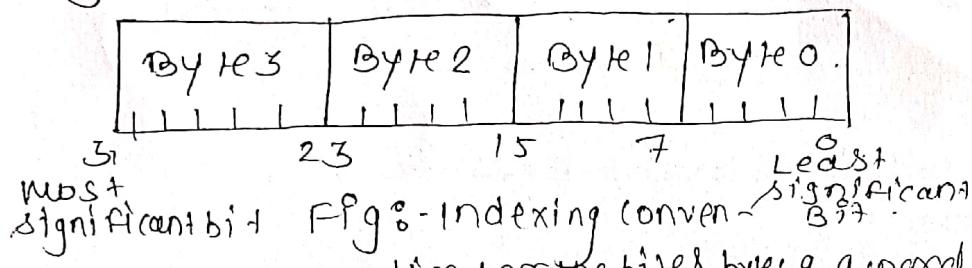
The CPU also has a standard word size for the data it processes. Word size is typically a multiple of 8, common CPU word sizes being 8, 16, 32 & 64 bits.

bits	Name	Illustration	Typical uses
1	Bit		Status flag. Logic variable
8	Byte		smallest addressable memory item. Binary coded decimal digit pair.
16	Halfword		short fixed-point no. Shows address (offset). Short instruction.
32	Word		Fixed / floating point no. memory Add: instrn.
64	Double word		Long instruction. Double-precision.

Fig!- Some information formats of the Motorola 680X0 microprocessor series.

→ Storage Order.

A small but important aspect of data representation is the way in which the bits of a word is indexed



→ we will usually follow the convention illustrated in above

where the eight-most bit is assigned the index 0 & the bits are labeled in increasing order from right to left.

Advantage of this convention is that when the word is interpret as an unsigned binary integer, the low-order indexes correspond to the numerically less significant bits & the high-order indexes correspond to the numerically more significant bits.

2. Fixed point numbers.

→ In selecting a no. representation to be used in a computer, following factors should be taken into account:

- The no. types to be represented : eg. int / real no.
- The range of values likely to be encountered.
- The precision of the numbers, which affects to the maximum accuracy of the representation.
- The cost of the hardware required to store & process the number.

→ Binary numbers :

- The fixed-point format is derived from the ordinary representation of a number as a sequence of digits separated by a decimal point.

The digits to the left of the decimal point represent an integer, the digits to the right represent a fraction. This is positional notation in which each digit has a fixed weight according to its position relative to decimal point.

If i^{th} then i^{th} digit to the left (right) of decimal point has weight (10^{-i}):

Thus five digit decimal number 192.73 is equivalent to

$$1 \times 10^2 + 9 \times 10^1 + 2 \times 10^0 + 7 \times 10^{-1} + 3 \times 10^{-2}$$

The most fundamental no. representation used in computer employs a base-2 positional notation.

A binary word of the form:

$$b_N \dots b_3 b_2 b_1 b_0, b_1, b_2, b_3, b_4, \dots b_m$$

Represents the number

$$\sum_{i=M}^N b_i 2^i$$

→ Signed Number.

Decimal Representation	Binary code			
	sign magnitude	ones complement	twos complement	
+7	0111	0111	0111	
+6	0110	0110	0110	
+5	0101	0101	0101	
+4	0100	0100	0100	
+3	0011	0011	0011	
+2	0010	0010	0010	
+1	0001	0010	0010	
+0	0000	0000	0000	
-0	1000	1111	0000	
-1	1001	1110	1111	
-2	1010	1101	1110	
-3	1011	1100	1101	
-4	1100	1011	1100	
-5	1101	1010	1011	
-6	1110	1001	1010	
-7	1111	1000	1001	

Table 8:-
Comparisons of
three 4 bit
codes for signed
binary number.

Table illustrates how integers are represented using all the codes when $n = 4$. These codes are all referred to as binary codes to distinguish them from the so-called decimal codes.

→ Decimal Numbers.

To converting an unsigned binary number $x_{n-1}x_{n-2}\dots x_0$ decimal requires a polynomial of the form.

$$\sum_{i=0}^{n-1} x_i \cdot 2^{k+i} \text{ to be evaluated.}$$

Eg: consider the addition $5+9=14$ using Excess - three code.

$$\begin{array}{r}
 & \begin{array}{c} 1000 = 5 \\ + 1100 = 9 \\ \hline \end{array} \\
 \text{carry, } & \leftarrow \begin{array}{c} 0100. \quad \text{Binary sum} \\ + 0011 \quad \text{correction} \\ \hline 0111 = 4 \quad \text{Excess-three sum.} \end{array}
 \end{array}$$

→ Hexadecimal Numbers.

- one / two other numerical codes are encountered in the design or use of computers. A particular importance is hexadecimal. The unsigned hexadecimal integer 2FAOC has the interpretation

$$\begin{aligned}
 & 2 \times 16^4 + F \times 16^3 + A \times 16^2 + 0 \times 16^1 + C \times 16^0 \\
 & = 2 \times 65,536 + 15 \times 4,096 + 10 \times 256 + 0 \times 16 + 12 \times 1 \\
 & = 195,084
 \end{aligned}$$

$$\text{Hence } 2FAOC_{16} = 195,084_{10}.$$

- Hexadecimal code is useful for representing long binary numbers. A consequence of the fact that the base is a power of 2.

A hexadecimal number is converted to binary simply by replacing each hex digit by the equivalent 4 bit binary form.

Eg) we can convert 2FAOC₁₆ to binary by replacing the first digit 2 by 0010, the second digit F by 1111, the third digit A by 1010, the fourth digit O by 1100.



3. Floating point Number.

The range of numbers that can be represented by a fixed-point number mode is insufficient for many applications, particularly scientific computations where very large & very small numbers are encountered. e.g. It is easier to write a quintillion as 1.0×10^{18} .

i) Basic formats.

Three numbers are associated with a floating point number a mantissa M, an exponent E and a base B. The mantissa M is also referred to as the significand / fraction in the literature.

These three components together represent the real number $M \times B^E$.

For example in its implementation the mantissa & exponent are encoded as a fixed-point number with a base β usually $2/10$.

A floating point number is therefore stored as a word (M,E) consisting of a pair of signed fixed-point numbers, a mantissa M, which is usually a fraction (an integer) & an exponent E, which is an integer.

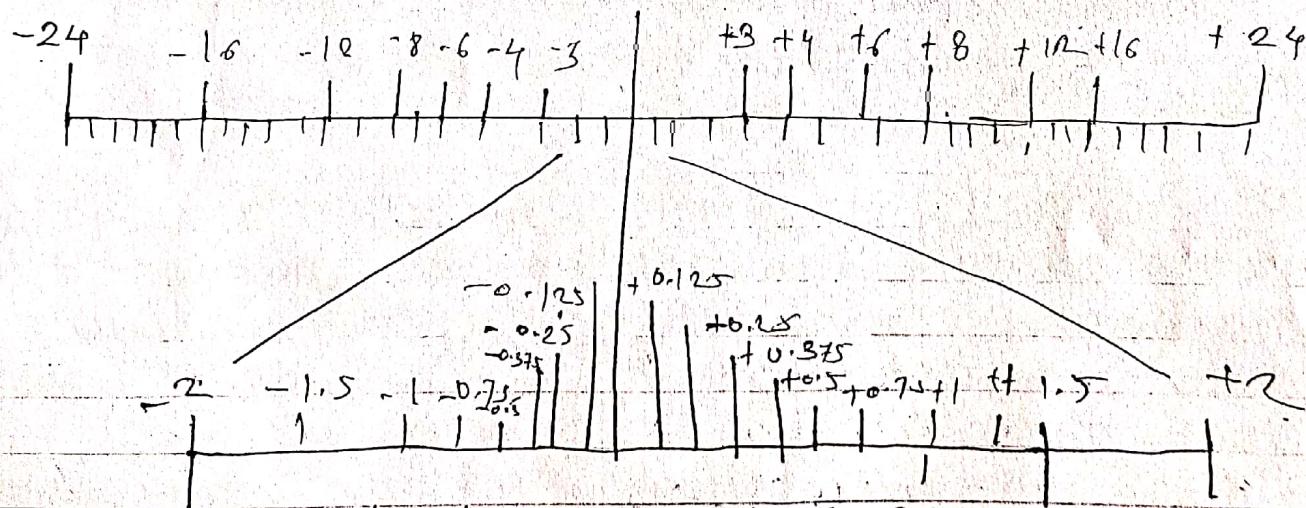


Fig. The Real numbers representable by a hypothetical 6-bit floating-point format.

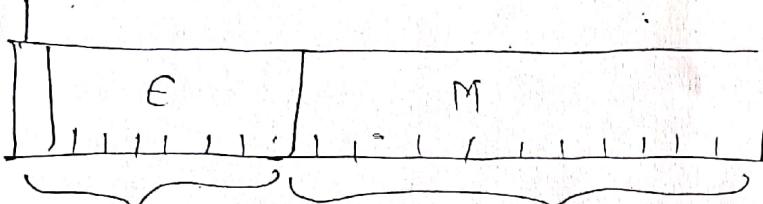
ii) Normalization and Biasing

- Floating point representation is redundant in the sense that the same number can be represented in more than one way.
- It is generally desirable to have a unique / normal form for each representable number in floating point systems.
- consider the common case where the mantissa is a sign magnitude fraction & base $q = 2$ is used. The mantissa is said to be normalized if the digit to the right of the radix point is not zero. i.e. no leading zeros appear in the magnitude part of the number.
- Normalization restricts the magnitude of a fractional binary mantissa to the range $1/2 \leq |M| < 1$.
- If k bits are allowed for the exponent including its sign, then 2^k exponent bit patterns are available to represent signed integers.
- The quantity K is called bias & an exponent encoded in this way called a biased exponent / characteristic.

iii) IEEE 754

- Until the 1980's floating point number formats varied from one computer family to the next, making it difficult to transport programs between different computers without encountering small but significant differences in such areas as round-off errors.
- To deal with this problem, the IEEE sponsored standard format for 32-bit and larger floating point numbers known as the IEEE 754 standard, which has been widely adopted by computer manufacturers.

SIGNS



($E_{\text{Excess-127}}$
binary integers))

C fraction part of sign
magnitude binary significand
(with hidden bit).

Fig. IEEE 754 Standard 32-bit floating-point number format.

The number N represented by a 32-bit IEEE standard floating point number has the following set of interpretations:

If $E = 255 \wedge M \neq 0$, then $N = \text{NaN}$.

If $E = 255 \wedge M = 0$, then $N = (-1)^S \infty$.

If $0 \leq E < 255$, then $N = (-1)^S 2^{E-127} (1.M)$.

If $E = 0 \wedge M \neq 0$, then $N = (-1)^S 2^{-126} (0.M)$.

If $E = 0 \wedge M = 0$, then $N = (-1)^S 0$.

Fixed point Arithmetic

Addition and Subtraction

Add & subtract instructions for fixed-point binary numbers are found in the instruction set of every computer.

In smaller machines such as microcontrollers they are the only available arithmetic instructions.

→ Basic Adder

The sum of z_i, c_i & 2 1-bit no. x_i & y_i can be expressed by the half-adder logic equations.

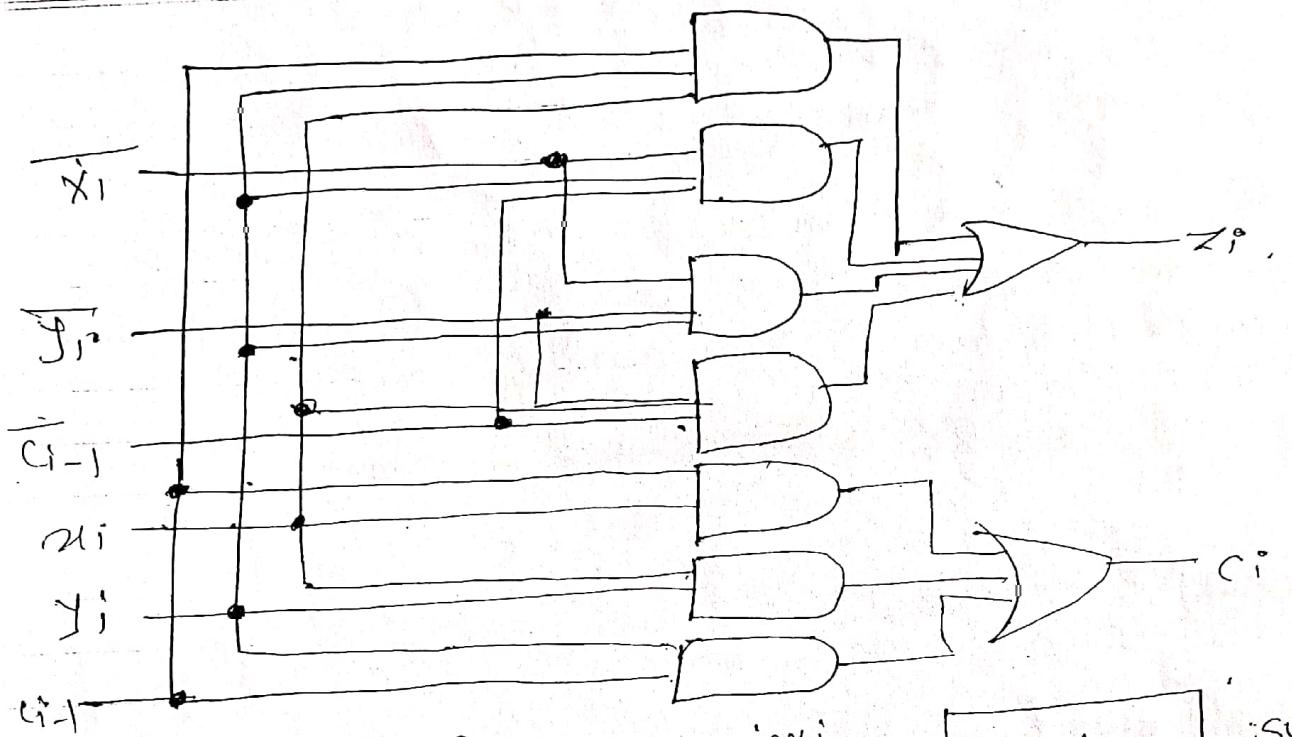
$$z_i = x_i \oplus y_i$$

$$c_i = x_i y_i$$

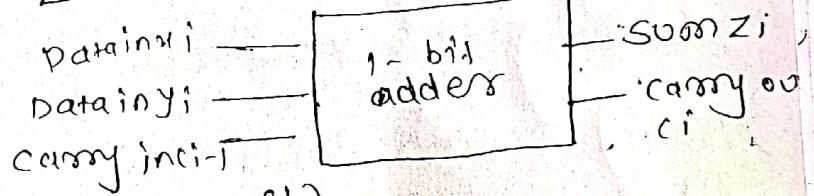
z_i = sum bit

c_i = carry-out bit

\oplus = exclusive-OR



(a).



(b).

Fig: A 1-bit (full) adder : (a) two-level AND-OR logic ckt & (b) Symbol.

→ Subtractors.

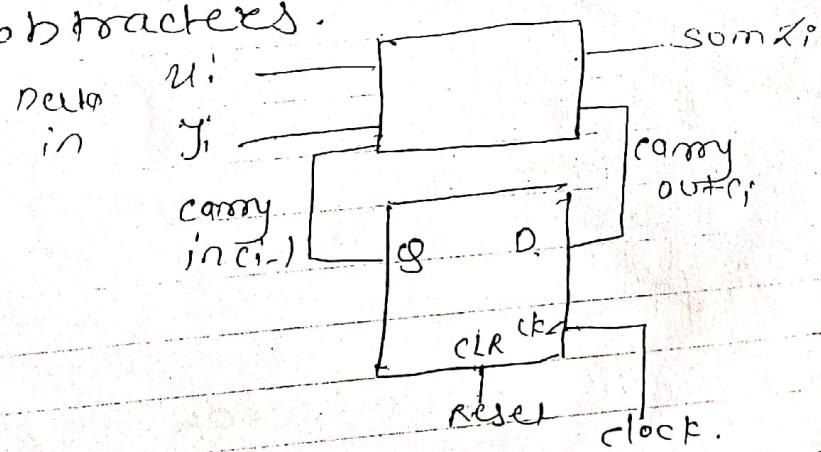


figure: A serial binary Adder

- Adding $-X$ to Y is equivalent to subtracting X from Y .
The ability to add negative numbers implies the capability

- Subtraction is relatively simple with two's-complement code because negation (changing X to $-X$) is very easy to implement.

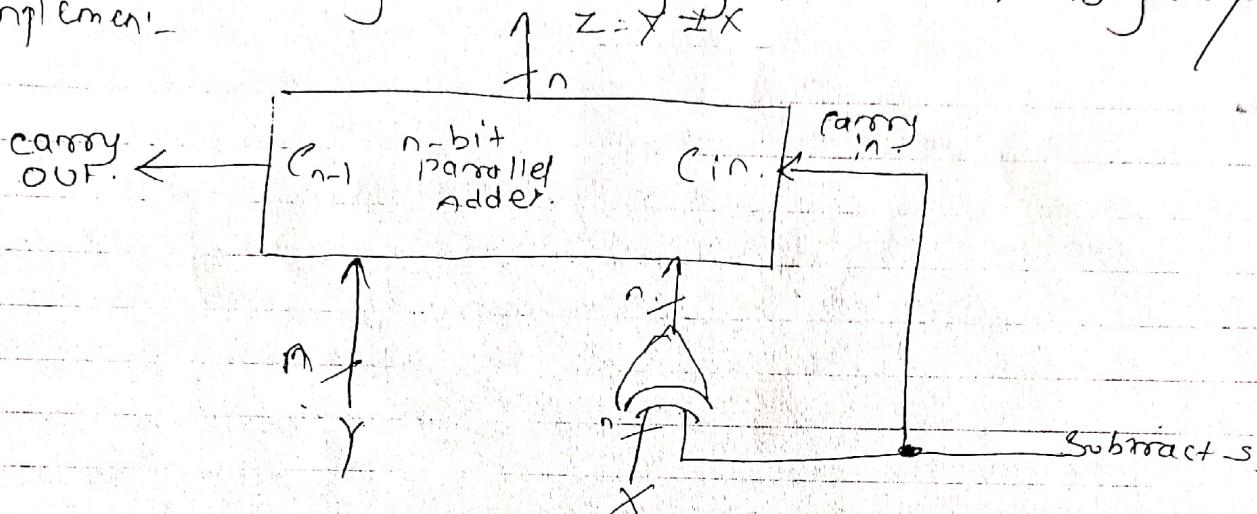


Fig: An n -bit two's-complement adder-subtractor

- 1) overflow : when the result of an arithmetic operation exceeds the standard word size n , overflow occurs
- overflow is indicated by a flag bit V in operations involving signed numbers: this flag is found in CPU status register.
- If we reinterpret the numbers in the preceding example as two's-complement rather than as unsigned, then $X = 11101011$ denotes -21_{10} , while $Y = 00101010$ denotes $+42_{10}$.
- The result Z computed in $Z = X + Y = 11101011 + 00101010 =$ now denotes $+21_{10}$ & the fact that $C_{n-1} = 1$ does not indicate overflow.
- In fact, we can never have overflow on adding a positive to a negative number.
- Overflow in modulo- 2^n two's-complement addition can only result from adding two positive numbers or two negative numbers.

iii) High-speed adder:

- The general strategy for designing fast adder is to reduce the time required to form carry signals.
- One approach is to compute the input carry needed by stage i directly from carrylike signals obtained from all the preceding stages $j=1, 2, \dots, i-1$ rather than wait for normal carries to propagate slowly from stage to stage. Adders that use this principle are called "carry lookahead adders".

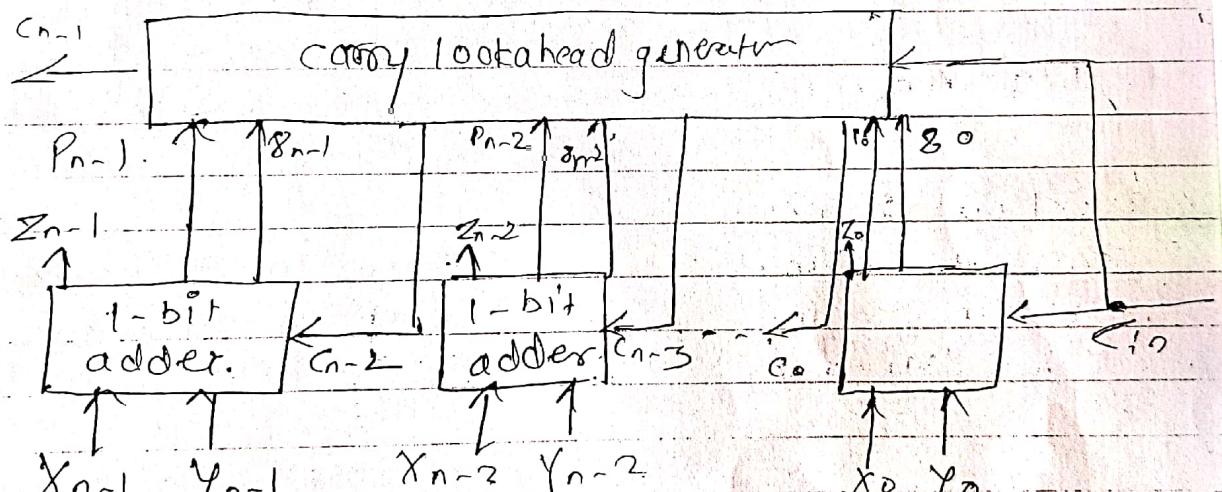


Fig. Overall structure of carry lookahead adder

Figure shows the general form of a carry lookahead adder circuit designed in this way.

We can further simplify the design by noting that the sum equation for stage i

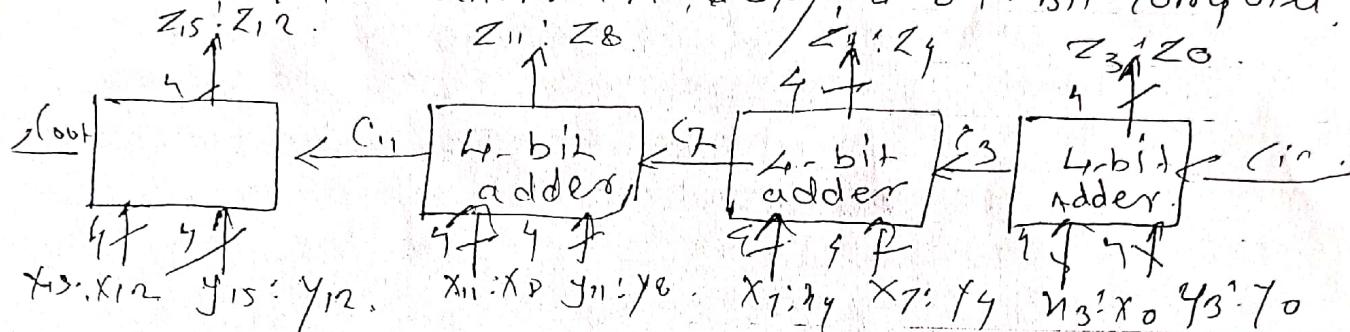
$$Z_i = x_i \oplus y_i \oplus c_{i-1}$$

is equivalent to

$$Z_i = p_i \oplus g_i \oplus c_{i-1}$$

V) Adder Expansion

The methods of handling carry signals in the two main combinational adder designs considered so far, namely Ripple carry propagation & carry-lookahead can be extended to large adders of the kind needed to execute add instructions in, say, a 64-bit computer.



5 Fixed point Multiplication

Fixed-point multiplication requires substantially more hardware than fixed-point addition and, as a result, is not included in the instruction set of some smaller processors. Multiplication is usually implemented by some form of repeated addition.

A simple but slow method to compute $x \cdot y$ is to add the multiplicand y to itself x times, where x is the multiplier. Often multiplication is implemented by multiplying y to x k -bits at a time and adding the resulting terms.

$$\begin{array}{r}
 1010 \text{ multiplier } y \\
 1101 \text{ multiplicand } x = 2^3 + 2^2 + 2^1 + 2^0 \\
 \hline
 1010 \text{ no } y \\
 0000 \cdot 2^3 y \\
 1010 \cdot 2^2 y \\
 1010 \cdot 2^1 y \\
 \hline
 10000010 \cdot \text{ product } p = \sum_{j=0}^3 2^j y
 \end{array}$$

Fig. typical pencil & paper method for multiplication of signed binary numbers.

i) Two's complement Multiplier

- The multiplication of two's complement numbers poses some difficulties in the case of negative operands.
- A conceptually simple approach to two's complement multiplication is to negate all negative operands at the beginning, perform unsigned multiplication on the resulting numbers & then negate the result if necessary. Two's complement negation for an integer $X = X_{n-1} X_{n-2} X_{n-3} \dots X_1 X_0$ is specified by

$$-X = \overline{X_{n-1}} \ \overline{X_{n-2}} \ \overline{X_{n-3}} \dots \ \overline{X_1} \ \overline{X_0} + 000 \dots 01$$

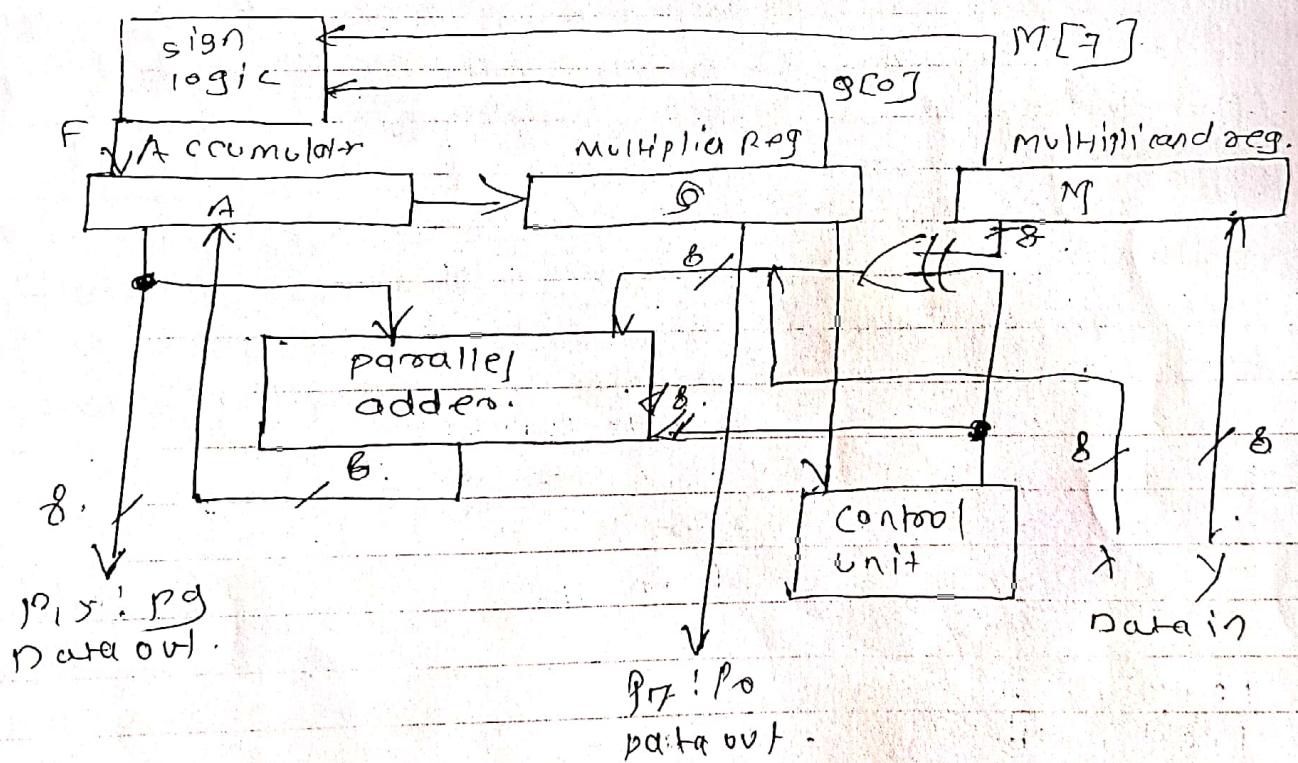


Fig : The datapath of the two's complement multiplier

ii) Booth's Algorithm

BoothMult

(in: INBUS; out: OUTBUS);

Registers A [7:0], M [7:0], Q[7:0], COUNT(2:0),
busINBUS [7:0], OUTBUS [7:0];

Begin:

INPUT:

SCAN:

A := 0, COUNT := 0

M := INBUS;

Q[7:0] := INBUS, Q[1] := 0;

If Q[1] Q[0] = 01 then A[7:0] := A[7:0],
M[7:0], go to TEST;

else if Q[1] Q[0] = 10 then A[7:0] := M[7:0],
A[7:0];

TEST : If COUNT = 7 then go to OUTPUT,

SHIFT : ACT := ACT, A[6:0], Q := A, Q[7:0],

INCREMENT : COUNT := COUNT + 1, go to SCAN;

OUTPUT : OUTBUS := A, Q[0] := 0;

OUTBUS := Q[7:0];

end BoothMult;

Fig: Description of an 8-bit multiplier implementing the basic Booth algorithm.

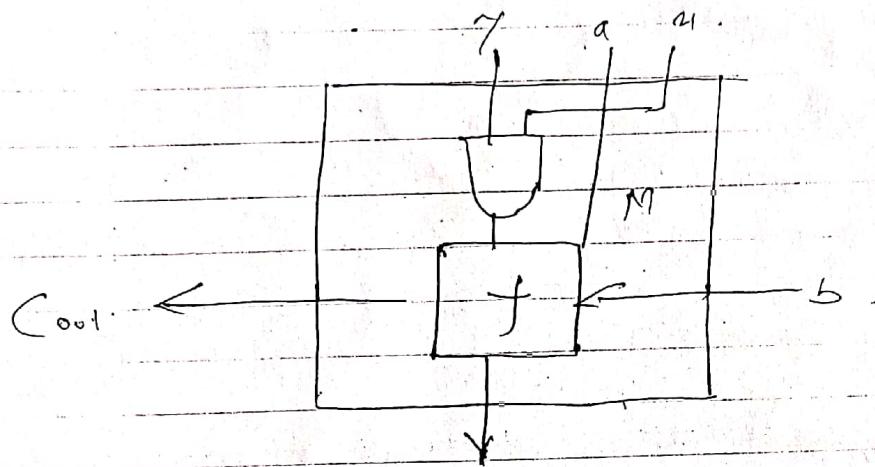
ii) combinational array multipliers.

Advances in VLSI technology have made it possible to build combinational circuits that perform $n \times n$ -bit multiplication for fairly large values of n .

Suppose that two binary numbers $X = x_{n-1}x_{n-2}\dots x_1x_0$ & $Y = y_{n-1}y_{n-2}\dots y_1y_0$ are to be multiplied. For simplicity, assume that X & Y are unsigned integers. The product $P = X \times Y$ can therefore be expressed as

$$P = \sum_{i=0}^{n-1} 2^i x_i \cdot Y$$





5 Fig: cell M for an unsigned array multiplier.

Fixed point DIVISION -

In fixed point division two numbers, a divisor V and dividend D, are given. The object is to compute a third or number Q, the quotient, such that $Q \times V$ equals or is very close to D.

$$\text{eq: } D = Q \times V + R.$$

where R = remainder, is required to be less than V , that is, $0 \leq R < V$. we can then write

$$D/V = Q + R/V$$

R/V = small quantity representing the error in using alone to represent D/V ; this error is zero if $R=0$.

i) Restoring and Non Restoring algorithm

The restoring and non-restoring division techniques will be extended to signed numbers in much the same way as multiplication.

- Nonrestoring division therefore requires n additions & subtractions, whereas restoring division requires an average of $\frac{3}{2}n$ additions & subtraction.

Rdivider : (in: INBUS; out: OUTBUS);
 Register S, A[n-1:0], M[n-1:0], Q[n-1:0],
 COUNT[logn]:0]; bus INBUS[n-1:0],
 OUTBUS(n-1:0) COUNT=0, S=0,
 EGIN ; A := INBUS; {Input the left half of the dividend}
 INPUT ; B := INBUS; {Input the right half of the divisor}
 BTRACT ; S.A := S.A - M; {s is the sign of the result}
 TEST ; If S = 0 Then
 begin Q[0] = 1;
 If COUNT = n-1 then go to CORRECTION; else
 begin COUNT := COUNT + 1, S.A.Q[n-1:n-1] := AG;
 end Else if s = 1
 begin Q[0] := 0;
 If COUNT = n-1 then go to CORRECTION; else
 S.A := S.A + M, go to TEST, end
 RRECTION ; If s = 1 then S.A := S.A + M.
 UTPUT ; OUTBUS := Q; {Q/R the quotient Q }
 OUTBUS := A; {Q/R the remainder R }
 end NRdivider;

Fig: Nonrestoring division Algo for unsigned integer.

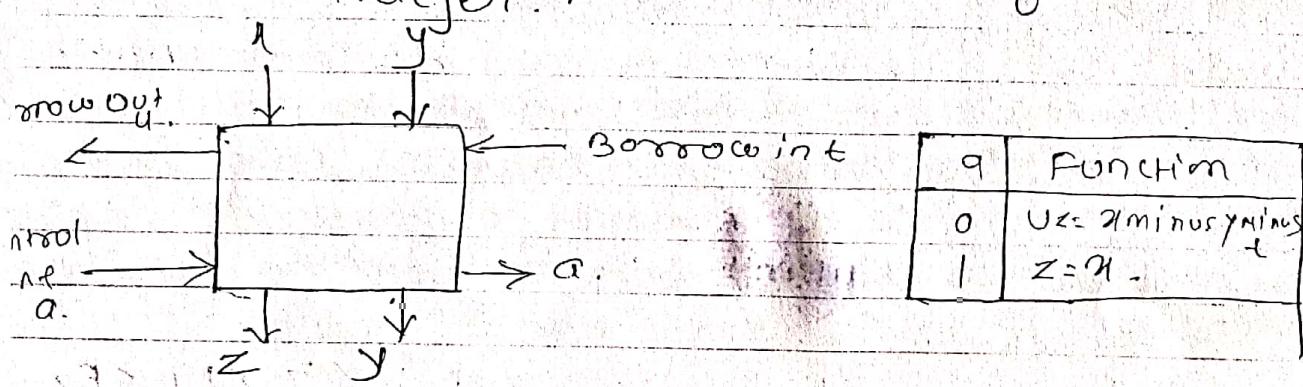


Fig: A cell D for array implementation of nonrestoring

i) combinational array divider:

combinational array circuits can be used for division as well as for multiplication.

Above fig. shows a cell \square suitable for implementing a version of the restoring division algorithm.

This cell is basically a full subtractor with b & u being the borrow-in and borrow-out bits, respectively.

The main output Z is controlled by input a . When Z is the difference bit defined by the arithmetic equation

$$Z = X \text{ Minus } Y \text{ minus } E$$

when $a=0$, $Z=X$. Thus the behavior of the cell is given by the logic equations.

$$Z = X \oplus \bar{a} (Y \oplus E)$$

$$u = \bar{a} Y + \bar{a} E + X E$$

iii) Division by repeated multiplication:

In system containing a high speed multiplier, division can be performed efficiently & at low cost using repeated multiplication. In each iteration a factor F_i is generated & used to multiply both the divisor V and the dividend D .

$$\therefore g = \frac{DXF_0 \times F_1 \times F_2 \times \dots}{V \times F_0 \times F_1 \times F_2 \times \dots}$$

F_i is chosen so that the sequence $V \times F_0 \times F_1 \times F_2 \dots$ converges rapidly towards one.

7. Floating point arithmetic

i) Basic operations

General formulas for floating point addition, subtraction, multiplication & division are given in below figure.

Addition

$$X+Y = (X_M 2^{x_e - Y_e} + Y_M) \times 2^{Y_e}$$

Subtraction

$$X-Y = (X_M 2^{x_e - Y_e} - Y_M) \times 2^{Y_e}$$

Multiplication

$$X \times Y = (X_M \times Y_M) \times 2^{Y_e + x_e}$$

Division

$$X/Y = (X_M / Y_M) \times 2^{x_e - Y_e}$$

Figure 8 :- The four basic arithmetic operations for floating-point numbers.

Floating point addition & subtraction have 3 main steps:

1. Compute $Y_e - X_e$, a fixed-point subtraction.
2. Shift X_M by $Y_e - X_e$ places to the right to form $X_M 2^{x_e - Y_e}$.
3. Compute $X_M 2^{x_e - Y_e} \pm Y_M$, a fixed-point addition/subtraction.

ii) Difficulties

- Several minor problems are associated with exponent biasing. If biased exponents are added / subtracted using fixed-point arithmetic in the course of a floating-point calculation, the resulting exponent is doubly biased & must be corrected by subtracting the bias.

Another problem arises from the all-0 representation usually required of zero.

A floating-point operation causes overflow or underflow if the result is too large or too small to be represented.

iii) Floating point units

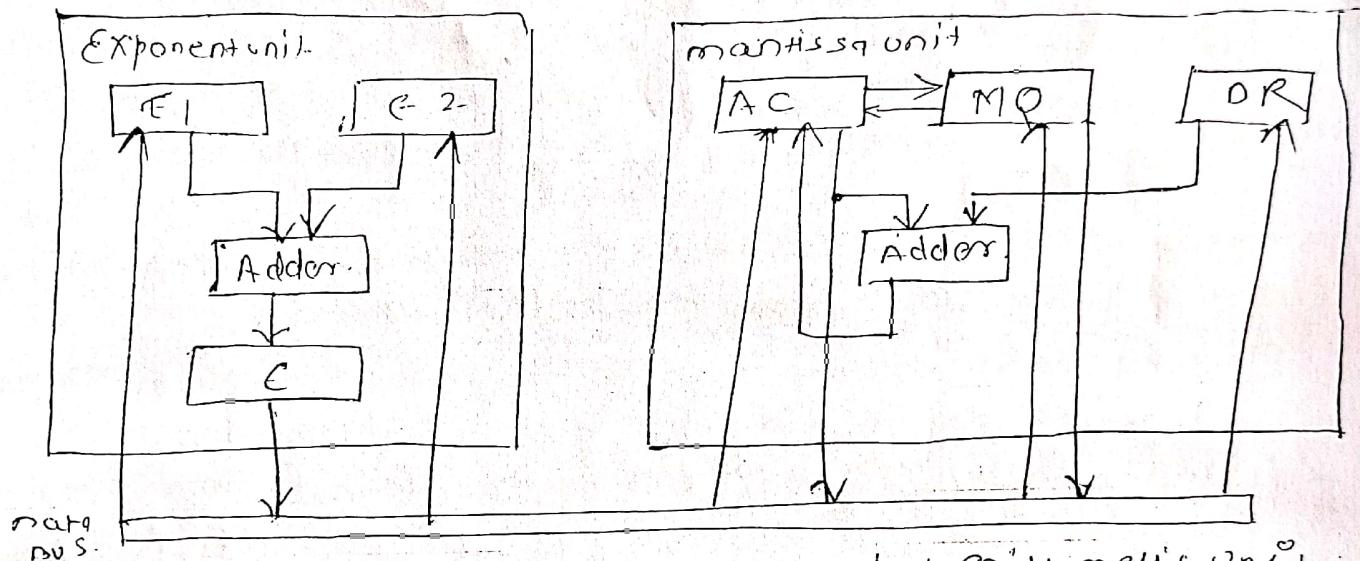


Fig : Datapath of a floating-point arithmetic unit

Floating-point arithmetic can be implemented by two loosely connected fixed-point datapath circuits, an exponent unit, and a mantissa unit.

iv) Addition.

The numbers have the 32-bit format of the IEEE standard 754. In this format each number N has a 23-bit fractional mantissa M with a hidden bit, an 8-bit exponent E is excess -127 code to a base 2^{128} . The value of N is therefore given by the formula

$$N = (-1)^E \cdot 2^{E-127} (1.M)$$

The numbers to be added in this instance are

$$\begin{array}{ll} X = 0 & 01111111 10000000000000000000 \\ V = 0 & 10000111 00101011010000000000000000 \end{array}$$