

Chapter 06:Computer Arithmetic

Hari K.C.

Department of Software Engineering
Gandaki College of Engineering and Science

Chapter 6: Arithmetic unit.

The arithmetic unit is a part of processor that executes arithmetic operations.

A microprocessor is capable of processing numeric data in one or more formats.

* Unsigned notation (unsigned representation of numbers)

⇒ In this notation, there is no any separate bit to represent the sign (positivity or negativity) of the number.

Signed representation of 8 bit data

e.g:-

1111

1 1001010

↓
sign magnitude

Sign = -7

But in unsigned representation, the whole 8 bit

unsigned
= 25 represent data.

The unsigned notations may be non-negative notations and 2's complement format.

✓ In non-negative notation, every number is treated as zero or a positive value. An n-bit number can have a value ranging from 0 (all bits are zero) to $2^n - 1$ (all bits are 1).

✓ In 2's complement format, both positive & negative numbers can be represented. An n-bit number can have range from -2^{n-1} to 2^{n-1}

$$\begin{array}{r} 127 \\ 2 \overline{(63)} \quad 1 \\ 2 \overline{(31)} \quad 1 \\ 2 \overline{(15)} \quad 1 \\ 2 \overline{(7)} \quad 1 \\ 2 \overline{(3)} \quad 1 \\ 2 \overline{(1)} \quad 1 \\ 2 \overline{(0)} \quad 1 \end{array}$$

negative numbers have 1 as the most significant bit and positive numbers (& zero) have 0 as leading bit.

<u>Binary representation</u>	<u>unsigned non-negative</u>	<u>unsigned 2's complement</u>
0000 0000	0	0
0000 0001	1	-1
0111 1111	127	-127
1000 0000	128	-128
1000 0001	129	-127
1	1	
0111 1110	254	-2
1111 1111	255	-1

$$\begin{aligned} 127 &= 0 \times 2^7 + 1 \times 2^6 + 1 \times 2^5 \\ &\quad + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 \\ &\quad + 1 \times 2^1 = 127. \end{aligned}$$

Addition & subtraction:

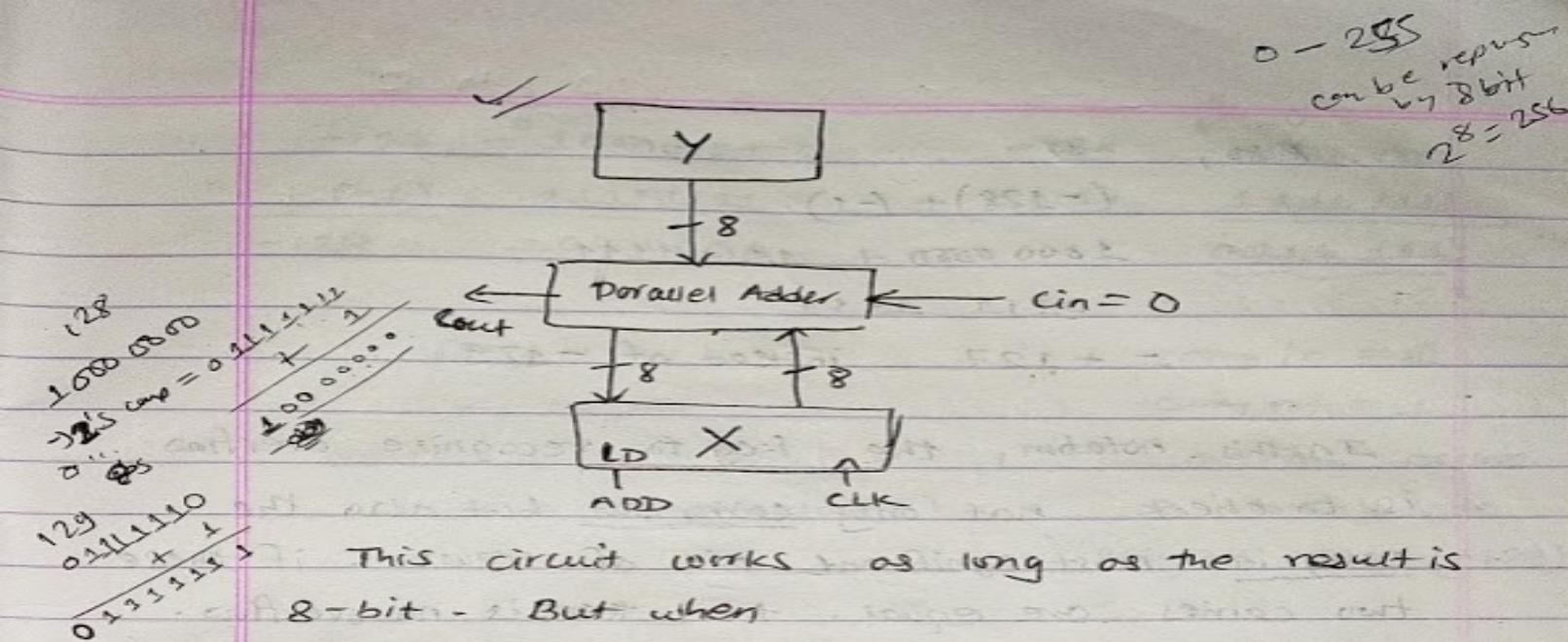
straightforward

Addition is realized by using a parallel adder.

$$ADD: x \leftarrow x + y$$

Here,

x & y are 8 bit registers



This circuit works as long as the result is 8-bit. But when

$$\begin{aligned}
 & 255 + 1 \\
 &= 1111\ 1111 + 0000\ 0001 \\
 &= 1\ 0000\ 0000 \Rightarrow 9\text{-bit value}
 \end{aligned}$$

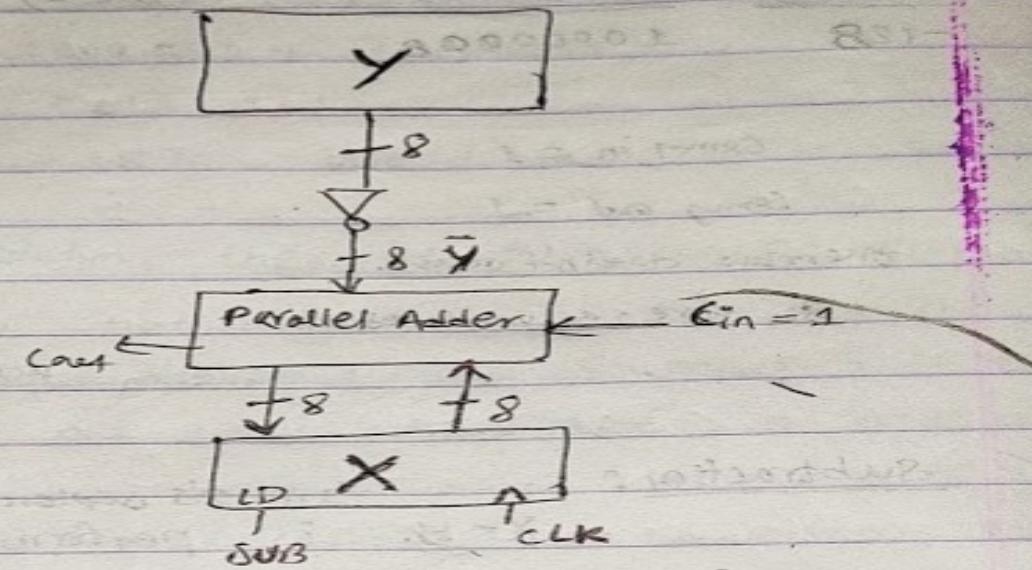
The 9-bit value cannot be stored in 8-bit register. The extra bit is cout which indicates arithmetic overflow. This bit can set an overflow flag.

In 2's complement notation, an overflow can occur at either side end of numeric range. At positive end, adding

$$\begin{aligned}
 & 127 + 1 \\
 & 0111\ 1111 + 0000\ 0001
 \end{aligned}$$

$= 10000000$. In 2's complement notation, this value is -128 not $+128$.

SUB: $X \leftarrow X - Y$



$$2 \rightarrow 0000\ 0010$$

$$\underline{(}-1\underline{)} \rightarrow \underline{+1111\ 1111}$$
$$\underline{\underline{1\ 0\ 0\ 0\ 0\ 0\ 0\ 1}}$$

$$255$$

$$- 254$$
$$\underline{\underline{0\ 0\ 0\ 0\ 0\ 1\ 0}}$$

$$1$$

$$- 2$$

$$\underline{+1111\ 1110}$$

$$\underline{\underline{0\ 0\ 0\ 0\ 0\ 0\ 1}}$$

$$01111\ 1111$$

multiplication:

$$\text{let } x = 27 \quad \times \quad 27$$

$$y = 253 \quad \underline{253}$$

$$\begin{array}{r} & 81 \\ 135 & \times \\ 54 & \times \\ \hline 6831 \end{array} \Rightarrow \begin{array}{l} \text{shift add} \\ \text{multiplication} \end{array}$$

shifts left

This normal calculator is calculated as

$$\underline{27}$$

$$\times \underline{253}$$

$$81$$

181 ← shifted one position right

$$27$$

$$\underline{253}$$

shifts right

$$81$$

81 ← shifted one position right

$$+ \underline{135}$$

$$1431 \leftarrow 1 \text{ is carried down, not added}$$

$$\underline{1431} \quad \leftarrow \text{shifted one position right}$$

$$\underline{54}$$

← 21 is carried down, not added

$$6831$$

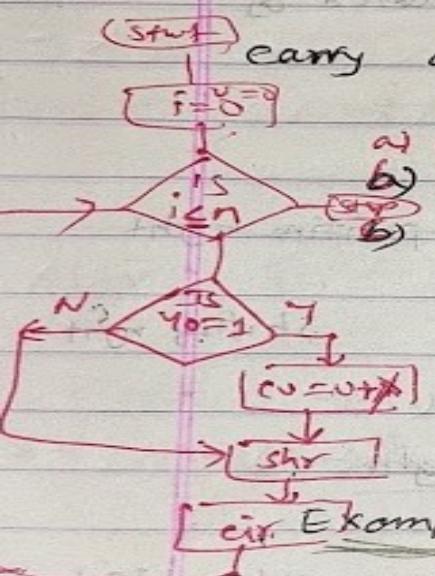
6831 one final shift

* Shift-add multiplication algorithm

⇒ Let X & Y be two n -bit registers (Input contents).

The result is stored in two n -bit register, U & V . U contains high order half of result & V contains low-order half.

C is 1-bit register used to store carry of an addition.



steps:

a) Start
Initialize $U=0, V=0$

b) Perform $i=1$ to n

- if $y_0 = 1$ Then $(C \cdot U) = V + X$
- linear shift right (CU, V)
- circular shift right (Y)

[cir. Example:-

$$\begin{array}{r} x \\ \times y \\ \hline 1101 * 1011 \end{array}$$

$$1101 * 1011 \quad \text{to } 801 \quad \text{carry to } 40$$

Function

$$\begin{array}{ccccccc} i & C & U & V & X & / \text{Comments} \\ \hline \end{array}$$

when $U=0$

if $y_0 = 1$

shr $((CU, V))$

cir (Y)

$$\begin{array}{ccccccc} & 0 & 1101 & 1000 & 1101 & / \text{w=(v+x)} \\ & 0 & 0110 & 1000 & 1101 & / \text{y}_0 = 1, \text{add} \\ & 1101 & & & & \end{array}$$

$$\begin{array}{r}
 0110 \\
 +1101 \\
 \hline
 10011
 \end{array}$$

if $y_0 = 0$, not
perform
 $CU = U + X$

since

if $y_0 = 1$ 2 1 0011

$\text{shy}(CUV)$

$\text{cir}(Y)$

Since $y_0 = 0$

new $y_0 = 0$

if $y_0 = 1$ 3

$\text{shy}(CUV)$

$\text{cir}(X)$

since

if $y_0 = 1$ 4 1 0001

$\text{shy}(CUV)$

$\text{cir}(Y)$

DONE

$y_0 = 1$ add
 $CU = U + X$

$$\begin{array}{r}
 1110 \\
 \hline
 y_0 = 0
 \end{array}$$

$X + U - V \quad y_0 = 0$

$$\underline{0111}$$

(UUV) \rightarrow E's 8

(Y) \rightarrow $y_0 = 1$,
 $CU = U + X$

1011 original value

$\underline{1000} \quad \underline{1111}$

\rightarrow U \rightarrow Result = X43

$$\begin{array}{r}
 1000 \quad 1111 \\
 \hline
 U \quad V \quad \rightarrow \quad 10001111
 \end{array}$$

RTL code for this algorithm

$O = 0$ 1: $U \leftarrow 0, i \leftarrow n$

y_0 2: $CU \leftarrow U + X$

2: $i \leftarrow i - 1$

3: $\text{shy}(CUV), \text{cir}(Y)$

2'3: GOTO 2

23: FINISH $\leftarrow 1$

$$1101 \rightarrow 1's complement$$

$$\begin{array}{r} 0010 \\ + 1 \\ \hline 0011 \end{array} = -3.$$

- speed up multiplication

~~Above algorithm~~ doesn't work for

negative numbers.

For example: $\begin{array}{r} 1101 \\ \times 1011 \\ \hline \end{array}$ would represent -3×-5 respectively in 2's complement notation

Booth algorithm directly works on two's

complement numbers, eliminating conversion b/w positive & negative numbers.

$[U=0, Y=0]$

Steps:-

i) Start

1-bit register

ii) Initialise $U=0, Y_{-1}=0$

iii) Perform for $i=1$ to n

a) if start of a string of 1's

in Y then $U = U - X$

i.e Y_0, Y_{-1}

if end of a string of 1's

in Y then $U = U + X$

i.e Y_0, Y_{-1}

b) Arithmetic shift right (UV)

c) circular shift right (Y) AND

copy Y_0 to Y_{-1}

$$\begin{array}{r} 0000 \\ - 1101 \\ \hline 0011 \end{array}$$

$$\begin{array}{r} 0000 \\ \times 0010 \\ \hline 0000 \end{array}$$

$$-8x - 5 = 15$$

$$\begin{array}{r} x = -3 \\ y = -5 \\ = 1101 \text{ of } 1011 \end{array}$$

$$2x + 5 = 15$$

Trace of Booth's algorithm

Function

i
loop

U=0

$y_{-1} = 0$

if start of string

1

if end of string

ashr(UV)

cir(Y), $y_{-1} \leftarrow y_0$

if start of string

2

if end of string

ashr(UV)

cir(Y), $y_{-1} \leftarrow y_0$

if start of string

3

if end of string

ashr(UV)

cir(Y), $y_{-1} \leftarrow y_0$

if start of string

4

if end of string

ashr(UV)

cir(Y), $y_{-1} \leftarrow y_0$

DONE

$$\begin{array}{r} U = 0 \\ V = 1101 \\ \hline 0000 & 0000 & 1011 & 0 \end{array}$$

Comments

$U = U - X$
start of string

$$\begin{array}{r} 0011 \\ \times 1011 \\ \hline 0001 & 1000 \\ \hline 1101 & 1 \end{array}$$

still within string

$$\begin{array}{r} 0000 \\ \times 1101 \\ \hline 0000 & 1100 \\ \hline 1110 & 1 \end{array}$$

$U = U + X$
End of string

$$\begin{array}{r} 1101 \\ 1010 \\ \hline 0111 & 0 \end{array}$$

$U = U - X$
start of string

$$\begin{array}{r} 0001 \\ \times 1011 \\ \hline 0000 & 1111 \\ \hline 1011 & 1 \end{array}$$

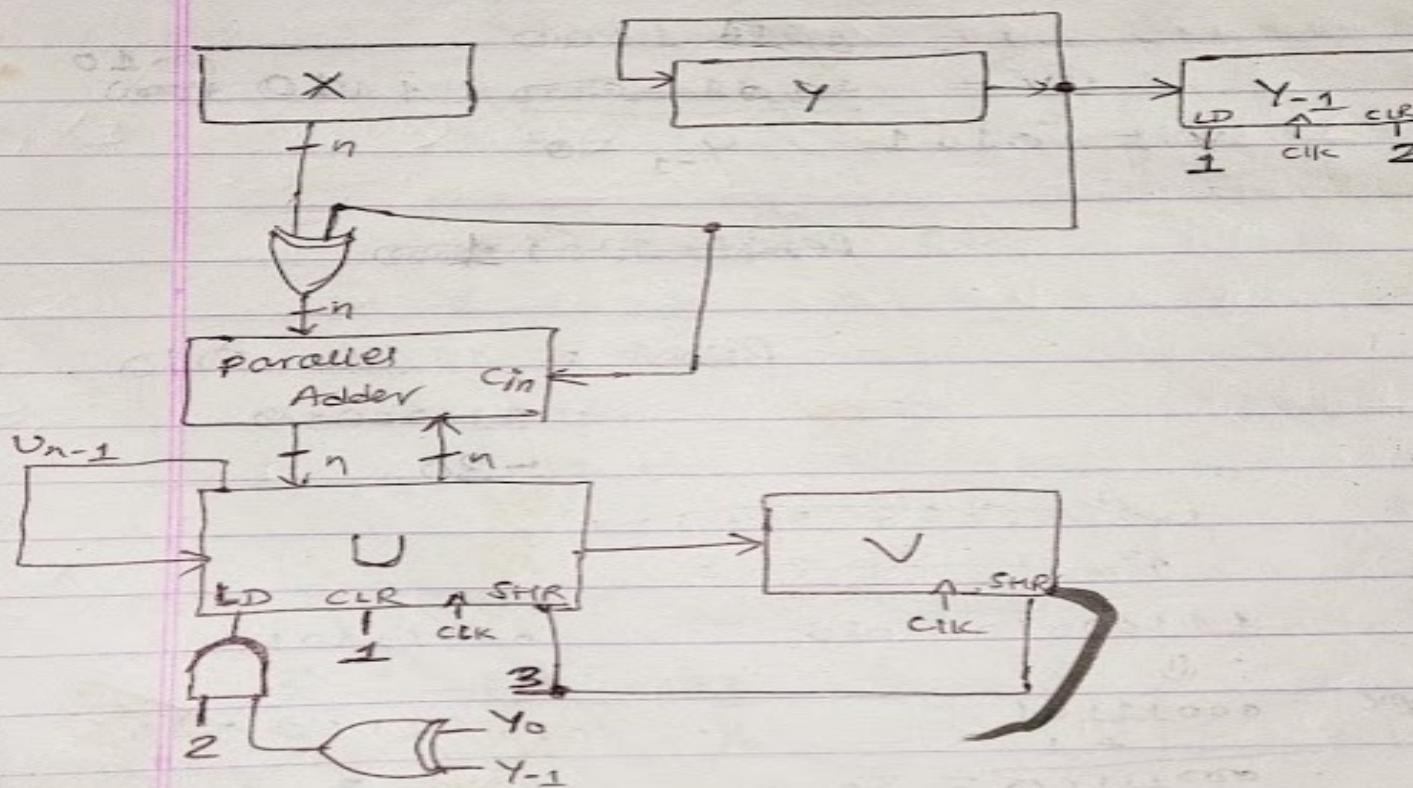
original value
Result = +15

$\text{ashr} \rightarrow$ arithmetic shift right
 $1100 \rightarrow 1120$
 $1 \rightarrow 1 \rightarrow 0 \rightarrow 1 \rightarrow 0$

RTL code for Booth algorithm

```

1: U<=0, Y-1<=0, i<=n
Y0 Y-1 2: U<=U+X+1 -> Y (Y0⊕Y-1) 2: U<=U+
Y0 Y-1 2: U<=U+X (X⊕Y) + Y
2: i <= i-1
3: ashr(UV), Cir(Y), Y-1 <= Y0
≥' 3: GOTO 2
≥3: FINISH <= 1
    
```



$$\begin{array}{r} 2 | 147 \\ 2 | 73 \\ 2 | 36 \\ \cancel{2} | 18 \\ 2 | 9 \\ 2 | 4 \\ 2 | 2 \end{array}$$

10010011

* Division

shift subtract division algorithm

The dividend is initially stored in UV. where U & V are n bits register. The divisor is stored in X register. The quotient in y register. The remainder will be stored in U. C is a 1-bit register to store extra bit shifted out of U.

steps:-

- i) Start
- ii) if $U \geq X$ then exit with overflow;
- iii) initialize $Y=0$ & $C=0$;
- iv) Perform for $i=1$ to n
 - f a) linear shift left ($C.UV$);
 - b) linear shift left (Y);
 - c) If $(C.U) \geq X$ Then { $Y_0 = 1$, $U = (U-X)$ }
cue movement true

= no. of bits

Example:

$$\begin{array}{r} 147 \div 13 \\ \hline 10010011 \quad \div \quad 1101 \\ \text{dividend} \quad \quad \quad \text{divisor} \\ U \quad V \quad \quad \quad X \end{array}$$

$$\begin{array}{r} 2 | 147 \\ \quad \quad \quad 7 \end{array}$$

$$\begin{array}{r} 10010 \\ - 1101 \\ \hline 0101 \end{array}$$

← shift
0

$$\begin{array}{r} 10010 \\ - 1101 \\ \hline 0101 \end{array}$$

Trace of shift subtract division algorithm

i	function	C	U	V	Y	comments
Initial		0	1001	0111	0000	
	if $U \geq X$			left		$U < X$, No exit
1	shl(CUV)	1	0010	0110	0000	
	shl(Y)				0000	
	if $CU \geq X$	0101		0010	0001	
2	shl(CUV)	0	1010	1100	0010	
	shl(Y)			0010	1000	
	if $CU \geq X$		(no operation)			$CU < X$ $01010 < 1101$
3	shl(CUV)	1	0101	1000	0100	
	shl(Y)			0100	1010	
	if $CU \geq X$	1000		0101		$10101 \geq 1101$ $U = 10101 - 1101$ $U = 1000$
4	shl(CUV)	1	0001	0000	1010	
	shl(Y)			1010		$10001 \geq 1101$
	if $CU \geq X$	0100		1011	Quotient	
			↓			
			remainder			

$$\begin{array}{r} 10001 \\ - 1101 \\ \hline 0100 \end{array}$$

RTL code for shift subtract division algorithm

Here X , U , V and Y are n -bit values

C & ~~overflow~~^{OVERFLOW} are 1 bit values

$x = 1, 2, 3, 4$ are sequence of states

G1 : FINISH $\leftarrow 1$, OVERFLOW $\leftarrow 1$

2 : $Y \leftarrow 0$, $C \leftarrow 0$, OVERFLOW $\leftarrow 0$, $i \leftarrow n$

3 : SHL(CUV), SHL(Y), $i \leftarrow i-1$.

(C+G) 4 : $Y_0 \leftarrow 1$, $U \leftarrow U + \bar{X} + 1$

Z4 : GOTO 3

Z4 : FINISH $\leftarrow 1$

The RTL code for the trace of the operation of $147 \div 13$

$U \rightarrow 1001$

$V \rightarrow 0011$

$X \rightarrow 1101$

$n = 4$

when $i = 0$ then

Z becomes 1.

* Signed notations:

In signed notations, most significant bit represent sign of the numbers (0 for positive number & 1 for negative number) and others bit represent the magnitude.

$$+3 = \underline{0} \text{ 0011}$$

$$-3 = \underline{1} \text{ 0011}$$

* Addition & subtraction of signed magnitude numbers.

8 4 2 1

$$\begin{array}{r} 0101 \\ +1x1 \\ \hline 1010 \\ = 10+1=5 \end{array}$$

* BCD:

Binary coded decimal.

It is the most popular format to represent decimal data.

In BCD, each 4 bits represent one decimal digit. $0001 \rightarrow 1$, $1001 \rightarrow 9$. From values 0 to 9, their representations are same as in binary. 4 bit values above 9 are not used in BCD since they do not correspond to any decimal digit.

Multidigit decimal numbers are stored as multiple groups of 4 bits per digit.

Eg:- 11 23 (Decimal)

0001 0001 0010 0011 (BCD)

In BCD, $+14 = \frac{0}{\text{Signbit}} 00010100$

$-14 = \frac{1}{\text{Signbit}} 00010100$

Addition & Subtraction using BCD

As long as sum of two digits is no more than 9, a binary parallel adder will give correct result.

There are two situations when there will be invalid BCD result:-

i) $8 + 2 = 1000 + 0010 = 1010$

which is not a valid BCD o/p

ii) $8 + 9 = \frac{1}{\text{carry}} 0001$ $1000 + 1001 = 10001$
not a correct result

although it is a BCD value.

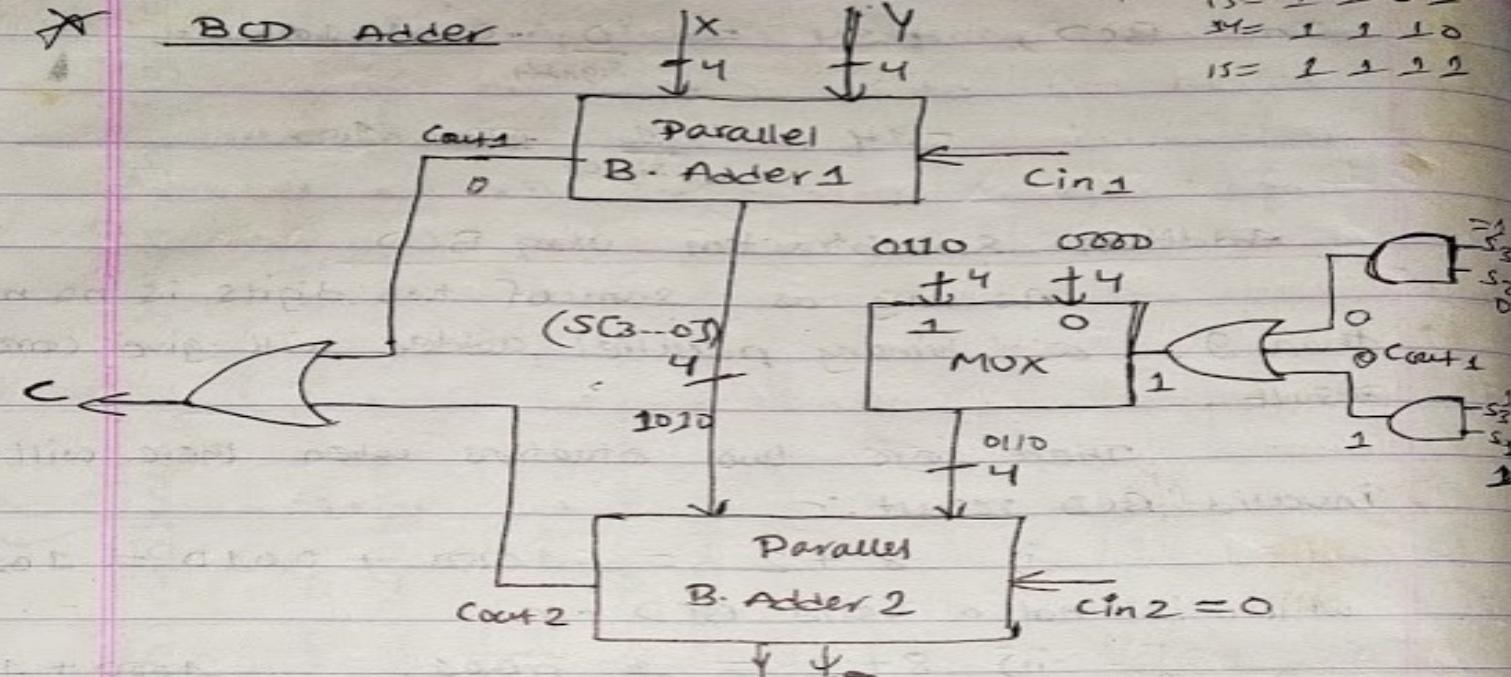
In both cases, adding 6 to the result generated by parallel binary adder produces correct BCD result.

~~$1010 + 0100 = 1010$~~

i) $1010 + 0100 = \underline{1} \underline{0000} \Rightarrow \text{BCD } (10)$
(6) 9 comes

ii) $0001 + 0110 = \underline{0111} \Rightarrow \text{BCD } (7)$

~~BCD~~ Adder



$S_3 \ S_2 \ S_1 \ S_0$
10 = 1 0 1 0
11 = 1 0 1 1
12 = 1 1 0 0
13 = 1 1 0 1
14 = 1 1 1 0
15 = 1 1 1 1

\Rightarrow In above figure $X \& Y$ (each of 4 bits) are added together. Then this result is added to either 0 (0000) or 6 (0110) to produce the correct BCD sum. If the result $S_3 S_2 S_1 S_0$ is not a valid BCD digit then either

$S_3 \wedge S_2 = 1$ or $S_3 \wedge S_1 = 1$. (This condition is true when invalid bits are from 1010 to 1111). In this case, multiplexer Select bit is set to 1, which causes 6 to be added to the result, correcting its decimal representation

The multiplexer select bit is also set to 1 if addition of X & Y generates carry out as when $X = 8$ & $Y = 9$. This also causes the circuit to add 6 to original sum.

If either parallel adder generates a carry out of 1, the carry out of BCD adder should also be 1.

* Disadvantages of using BCD notation

1) By representing numbers in decimal it is a waste of considerable amount of storage space

$$17 \text{ } 83 \rightarrow \begin{smallmatrix} 00010111 \\ 00001111 \end{smallmatrix} \text{ - BCD}$$

Since the no. of bits needed to store decimal number in binary code is greater than the no. of bits needed for its equivalent binary representation.

2) The circuit requires to perform decimal arithmetic are complex.

Despite of these limitations there are advantages of decimal representations.

Advantage of BCD

Computer inputs & outputs data are generated by people who use the decimal system. Some applications such as business data processing require small amount of arithmetic computations compared to the amount required for input & output for decimal data.

For this reason, some computers & all electronic calculators perform arithmetic operations directly with the decimal data (in binary code) and thus eliminate the need for conversion to binary & back to decimal.

Some computer systems have hardware for arithmetic calculating with both binary & decimal data.

* Specialized Arithmetic Hardware

➤ The different special arithmetic hardware are designed to speed up arithmetic computations. F0 processors were formerly used to speed up computations.

The coprocessor chip had specialized hardware that performed arithmetic calculations much more quickly than microprocessor. The coprocessor monitored instructions on system data bus. When the microprocessor fetched an instruction that the coprocessor could execute, the coprocessor sent a signal to microprocessor indicating that it would perform computation. The coprocessor then calculated the result & sent it to microprocessor. If the coprocessor was not present, the microprocessor executed the instruction in slow way.

* 3 techniques to improve speed of arithmetic calculation

- 1) Pipelining
- 2) Look up tables
- 3) Wallace tree multipliers

any Coca Cola manufacturing process.

1) Pipelining: Arithmetic Pipelining consists of number of stages. Firstly, data enters a first stage which performs some operations on data. Then results are passed to next stage which performs its operations. and so on until the final computations has been performed.)

For arithmetic pipeline, an individual stage might be an adder or a multiplier. A pipeline does not speed up individual computation. It improves performance by overlapping computations. For e.g:- In a 3 stage pipeline, the first stage process data set 1 while second stage process data set 2 & last stage process data set 3. The net effect is that results are output more quickly than in non-pipelined arithmetic.

For $i = 1$ TO 100 DO { $A[i] \leftarrow (B[i] \cdot C[i]) + D[i]$ }

Each operation of multiplication & addition takes 10 ns to complete.

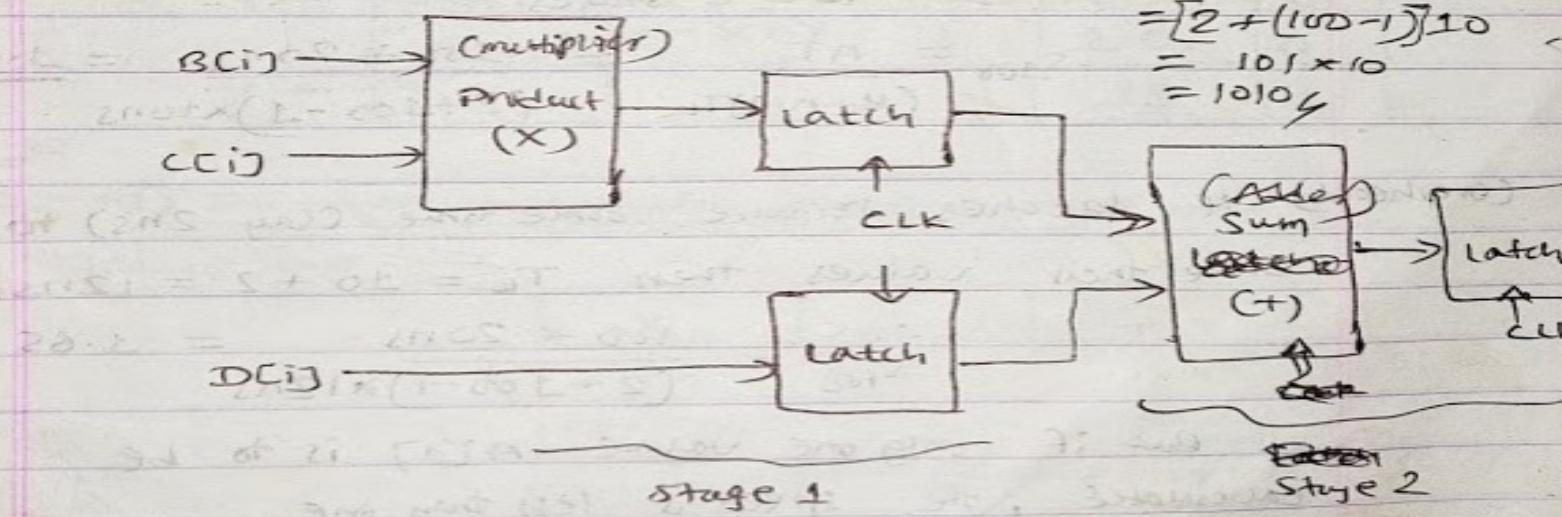
A non pipelined processor now takes 200ns to calculate $A[i]$ & total $(100 \times 200\text{ns}) = 20000\text{ns}$ to execute code.

$$\frac{B[1] * C[1] + D[1]}{\text{1st stage}} \xrightarrow{\text{10ns}} \text{D}[1] \quad B[2] * C[2] + D[2]$$

2nd stage

But a pipelined can perform $B[3]$ operation in two stages - The first stage perform multiplication & second perform addition.

In pipelining, during first 10ns, the first stage calculates $B[1] * C[1]$. In next 10ns, stage second stage calculates adds this value to $D[1]$ & store result in $A[1]$. At the same time when second stage is performing, stage 1 multiplies $B[2] * C[2]$ then after 10ns stage 1 performs $B[3] * C[3]$ & stage 2 calculates final value $A[2]$. So, the pipeline execute above code in 1010 ns. instead of 2000 ns



$$\begin{aligned}
 & (10 + (n-1))t_k \\
 &= [2 + (100-1)]10 \\
 &= 101 \times 10 \\
 &= 1010
 \end{aligned}$$

Stage 2

- some be with

The speedup factor of pipeline is given as

$$S_n = \frac{nT_1}{(k+n-1)T_k}$$



where S_n is amt. of time needed to process ' n ' pieces of data. k = stages. T_1 is amt. of time to process one piece of data. T_k is clock of k -stage pipeline.

In above example

$$T_1 = 20\text{ ns}$$

$$T_k = 10\text{ ns}$$

$$n = 100$$

$$k = 2 \text{ stages}$$

$$S_{100} = \frac{nT_1}{(k+n-1)T_k} = \frac{100 \times 20\text{ ns}}{(2+100-1) \times 10\text{ ns}} = 1.98$$

(overhead) if latches require some time (say 2ns) to store their values then $T_k = 10 + 2 = 12\text{ ns}$

$$S_{100} = \frac{100 \times 20\text{ ns}}{(2+100-1) \times 12\text{ ns}} = 1.65$$

But if only one value $A[17]$ is to be calculated, the speed is less than one

$$S_n = \frac{1 \times 20\text{ ns}}{(2+1-1) \times 12\text{ ns}} = 0.83 < 1$$

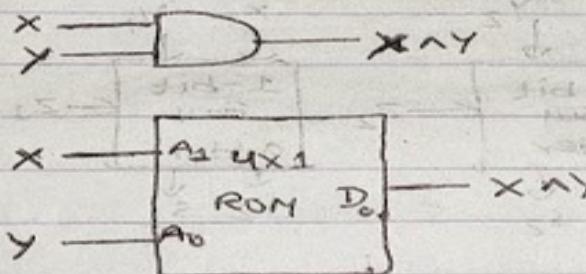
\Rightarrow slower than non-pipelined

2) LookUp Tables :-

A combinational circuit can be implemented by a ROM configured as look up table. The inputs to the combinational circuit serve as address inputs of ROM. The data outputs of ROM corresponds to output of combinational circuit.

Example:-

Let us consider 4×1 ROM programmed as two input AND gate

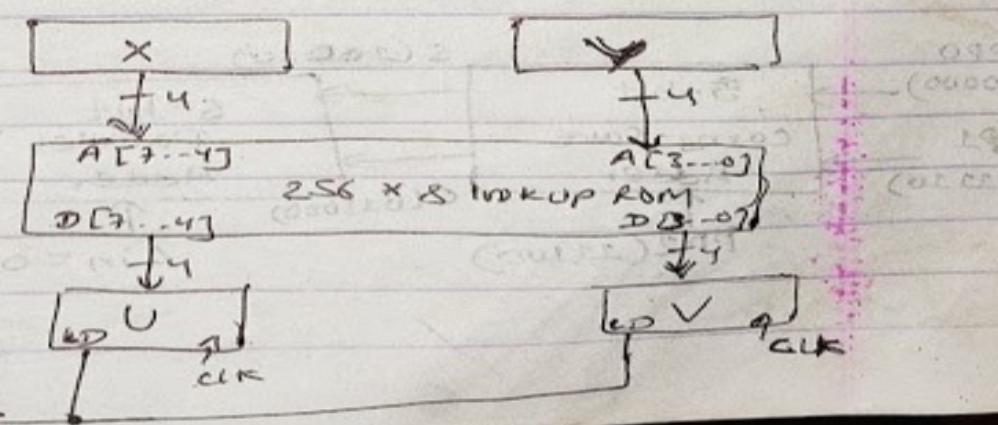


Address	Data
0 0	0
0 1	0
1 0	0
1 1	1

ie
doubles
the
speed

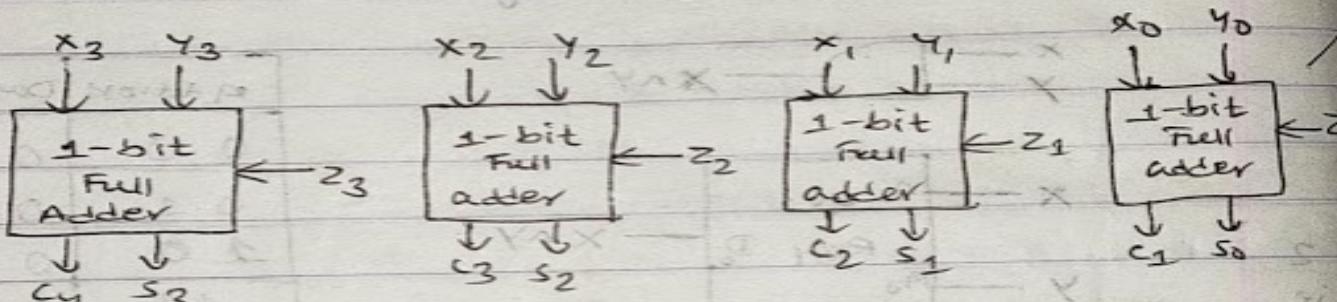
2) shift-add multiplication algorithm implemented by

look up ROM



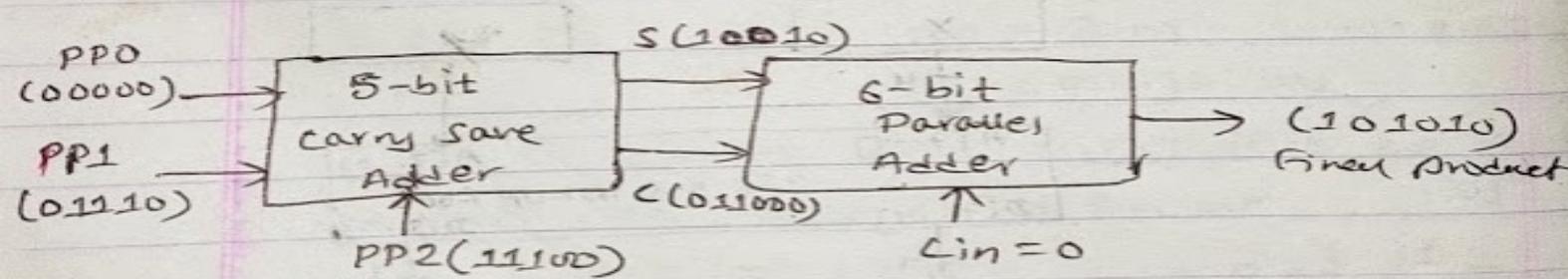
3) WALLACE TREES

A Wallace tree is a combinatorial circuit used to multiply two numbers. It requires more hardware but can produce a product in less time. It ~~uses~~ requires carry save adders & only one parallel adder. A carry-save adders can add 3 values simultaneously instead of two. It output both sum & a set of carry bits.



(fig : A carry save adder)

A 3×3 multiplier constructed with a carry save adder.



$\begin{array}{r} 0010 \rightarrow Z \\ 0111 \rightarrow X \\ + 1011 \rightarrow Y \\ \hline 1110 \end{array}$ sum

~~00110~~ carry
 ~~$c_0 = 1$~~ initial carry
 ~~$c_1 = 0$~~ carry 0
 ~~$c_2 = 1$~~ carry 1
 ~~$c_3 = 0$~~ carry 0

$x = 111$
 $y = 110$
 $z = 0010$
 $c_0 = 0$
 $c_1 = 1$
 $c_2 = 0$
 $c_3 = 1$

$\begin{array}{r} 1110 \\ 00110 \\ \hline 101010 \end{array}$ ← Final sum calculated.

Let us consider a example. $x = 0111$ (7)
 $y = 1011$ (11) and $z = 0010$ (2) this carry save adder produces

$s = 1110$ & $c = 00110$

$\begin{array}{r} 111 \\ 110 \\ 001 \\ 110 \\ \hline 101010 \end{array}$ ← Final sum

~~$11100 - PP2(2)$~~
 ~~$00000 - PP0(X)$~~
 ~~$01110 - PP1(Y)$~~

10010 sum
 01100 carry
 initial carry
 carry 0

$01100 \rightarrow C$
 $+ 010010 \rightarrow S$
 $101010 \rightarrow \text{final sum}$

* Floating point numbers

The floating point represent of a number has two parts. The first part represent a signed, fixed point number called mantissa. The second part designates the position of decimal (or binary) point and is called exponent. The fixed point mantissa may be fraction or an integer.

For e.g:- the decimal number
+ 6132.789 is represented in floating point
as :- fraction = 3 Exponent
 + 0.6132789 +04

The value of exponent indicates that actual position of decimal point is four positions to right of indicated decimal point in fraction.

This representation is equivalent to scientific notation $+ 0.6132789 \times 10^{+4}$
so, it has form
 $m \times r^e$

only mantissa m & exponent e are physically represented in register & r is the radix or base.

For e.g.: a binary number ~~+1001.110~~ is represented as fraction ~~x~~ binary exponent

$$\underline{0.1001110}, \quad 000100$$

The fraction has a 0 in leftmost position to denote positive (the binary point of fraction follows sign bit but is not shown). The exponent has equivalent binary number +4. The floating point number is equivalent to

$$m \times 2^e$$

$$= (0.1001110)_2 \times 2^{+4}$$

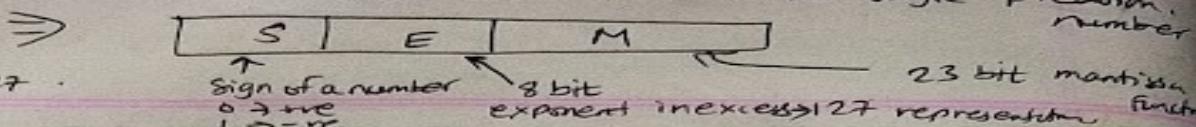
IEEE 754 floating point standard :

⇒ IEEE developed IEEE 754 floating point standard. This standard defines set formats & operation modes. IEEE 754 standard used in virtually all CPUs that have floating point capability.

IEEE 754 standard specifies two precision for floating point numbers.

Single precision numbers have 32 bits : 1 for the sign, 8 for exponent and 23 for the significant.

IEEE standard for floating point numbers for single precision.



The significand also include an implied 1 to the left of its radix point.

Double precision numbers use 64 bits:

1 for the sign, 11 for the exponent and 52 for significand

example (single precision)

+19.5

Sign = 0

Significand = 001110000000000000000000

Exponent = 100000011 (131)

(127 + 4)

bias value

Double precision

~~1001110~~

~~significand~~

$$+19.5 = 100111.1$$

$$= 1.00111 \times 2^4$$

not included.

2 | 19

2 | 9

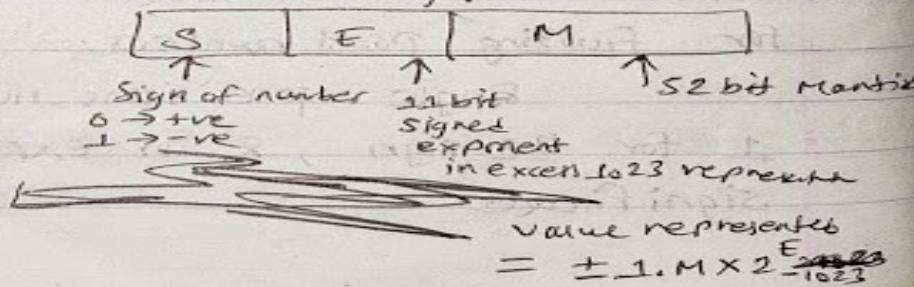
2 | 4

2 | 2

2 | 0

100111.10

IEEE standard for double precision floating point numbers.



THE END