## Chapter 3
### Arithmetic for Computers

\# এই chapter টা Calculative অর্থাৎ Calculation based

\# Overflow detect করার জন্য.

Arithmetic for Computers.

⇒ Operations on ~~Computer~~ integers.

- Addition & Substraction
- Multiplication & Division.
- Dealing with Overflow

⇒ Floating Point Real numbers
- Representation and operations.

> This chapter will deal with these things.

> \# এই chapter এ আমরা signed number নিয়ে deal করবো

## Addition & Substraction:

Integer Addition:

$7 + 6$

*For this example we are considering 6 bit arch.*

$$000 \quad 111 \longrightarrow \text{6 bit representation of +ve 7}$$
$$000 \quad 110 \Longrightarrow \text{6 bit representation of +ve 6}$$
$$\overline{001 \quad 101}$$

\# MSB bit 0 অর্থাৎ Positive

\# MSB bit 1 অর্থাৎ negative

\# দুইটা number addition বা substract করার সময়

wrong result আসলে overflow.

\# 2টা positive number add করলে result positive

আসার কথা, in some cases result negative

আসতে পারে ; কিন্তু তা তো হওয়ার কথা না. বিধায়

এখানে Overflow হয়েছে,

\# Similarly, 2টা negative number add করলে result

negative আসার কথা; but somehow result positive

আসলে Overflow.

• Integer Addition

• Integer Substraction.

একটা n bit system নিয়ে

কাজ করলে তার unsigned

number এর range

$-2^{(n-1)}$ to $+2^{(n-1)} -1$

# অস্বাভাবিক আমরা বিভিন্ন arithmetic operation করলে যে result আমরা পাই তা নাও হতে পারে ভুল result আসে। তার কারণ overflow.

## Integer Addition.

**\*\* 3টি কারণে Overflow হতে পারে,**

- If result out of range.

- Adding a positive number and a ~Positive ~~negative~~ number

  $\Big\{$ ↳ Overflow if $\boxed{\text{result sign}}$ is 1     → MSB bit

  ↳ একটি (+ve) এর সাথে একটি (+ve) যোগ করলে result (+) আসার কথা তার মানে result sign 0 আসার কথা.

- Adding two negative numbers.

  $\Big\{$ ↳ Overflow if result sign is 0

  ↳ দুইটি (-ve) number যোগ করলে result (-ve) আসার কথা; তার মানে result sign 1 হওয়ার কথা.

# যদি একটা positive number এর সাথে একটা negative number add করি তাহলে no overflow.

# প্রত্যেকটা calculation binary তে করছি, Signed binary number দিয়ে deal করছি।

---

Considering 4 bit architecture

7 + 6

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

↳ এখানে MSB bit টা 1 ; মানে 2 টা positive number add করার পর neg number পাচ্ছি। So, overflow.

→ Actually result টা 01101 অর্থাৎ তেরো কিন্তু 4 bit arch এর বাইরে 4 discard হয় তোমরা 1 so আর portion টা থাকায় overflow

---

# কত bit architecture follow করছি তার উপরে রাখতে হবে।

## Examples:

↳ নিচের example গুলোর জন্য 5 bit architecture consider করছি।

∴ range for 5 bit signed number:

$$-2^{(5-1)} \text{ to } 2^{(5-1)} - 1 = -16 \text{ to } 15.$$

@

$11 + (-1) = 10$

→ 10 is inside the range.

$$
\begin{array}{r}
0\,1\,0\,1\,1 \longrightarrow 11 \\
+\ 1\,1\,1\,1\,1 \longrightarrow -1 \\
\hline
\textcircled{1}\ \textcircled{0}\,1\,0\,1\,0
\end{array}
$$

↳ Updated MSB as 5 bit arch

Result 6 bit but system 5 bit so first এর 1 কাটা যাবে, so যদি 5 bit এর বেশি consider করবো।

Steps:
• আগের range বের করে নিতে হবে
• Then calculation করতে হবে,
• Check করতে হবে result out of range কিনা, যদি range এ আসে then check করে result sign টা। তারপর দিয়ে judgement হবে either it's overflow or not.

MSB bit এর result sign 0 so positive আসলে কথা positive এ আসা, so no overflow।

# Binary to Decimal Conversion জানতে হবে,

(P-7.0)

◫ Example 2    [Adding two positive numbers]

$$11 + 15 = 26$$

⇒ Initial decision: Overflow; as 26 is out of range.

$$01011 \longrightarrow 11$$
$$+ \quad 01111 \longrightarrow 15$$
$$\overline{\quad①1010\quad}$$

↳ MSB এর Result sign

$$\ominus\left\{(1 \times 2^4) + 1 \times 2^3 + (1 \times 2^1)\right\}$$
$$= -26$$

এ n ৬বে bcuz signed number

এখানে আমরা 2টা positive number add করছিলাম so result 3 positive আসবে করা; but calculation করার পর দেখা গেলো Result sign is 1; which indicates it's an negative number.

∴ So overflow confirmed.

◫ Example 3    [Adding two negative number]

$$-11 + (-15) = -26 \quad (Overflow)$$

$$10101$$
$$10001$$
$$\overline{①00110}$$

এই extra 1 টা discard হবে যাতে as 5 bit arch, বাকি portion নিয়ে check করবো,

result আসছে করা: negative। But result sign 0. S 0. Which indicates positive number.

∴ Overflow.

# Integer Substraction:

- 2 টা negative number এর মধ্যে substract করলে
  বা 2 টা positive number এর মধ্যে substract
  করলে, no overflow.

- Integer Substraction এ 3 টাতে overflow হতে
  পারে,

  ⇒ Out of range.

  ⇒ Substracting a positive value from a negative value

    - If result sign is 0 then overflow.

  ⇒ Substracting a negative value from a positive value.

    - If the result sign is 1 then overflow.

---

- 8 − (−4) = 12 তাহলে করা.

- −8 − (+4) = −12 তাহলে করা
  but আমরা positive so
  overflow.

# Examples:

☐  $3 - (-15) = 18$  [Substracting neg number from a pos number]

⟹ 18 is out of range; so, overflow.

$$\begin{array}{r} 0\ 0\ 0\ 1\ 1 \rightarrow 3 \\ +\ 0\ 1\ 1\ 1\ 1 \rightarrow 15 \\ \hline 1\ 0\ 0\ 1\ 0 \end{array}$$

$\left(\begin{array}{c} 3 - (-15) \\ 3 + 15 \end{array}\right.$

Result was supposed to be positive; but our result

sign is 1; which indicates negative number.

So, overflow.

☐  $-8 - (+10)$  [Substracting a positive number from a negative number]

$= -18$

⟹ $-18$ out of range so overflow.

$\left(\begin{array}{c} -8 - (+10) \\ -8 - 10 \\ -8 + (-10) \end{array}\right.$

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

↳ Result sign

As, result sign is 0; which was supposed

to be 1 in this case. So, overflow.

# Dealing With Overflow

**#** C language Overflow handle করতে পারে ।

    ↳ এর জন্য এটা যে MIPs এ " addu, addui, subu " instruction
     use করে।

**#** অন্যান্য language এ Overflow deal করার, এটা।

    exception raise করতে হয়,

    ↳ As তারা add, addi, sub instruction use করে।

**◩** So, how to deal or Overcome Overflow?

⇒ Follow the 4 steps stated below :

(i) If the current instruction have overflow then
     PC holds the address of that instruction.

(ii) Then PC stores the address of that instruction
     in the Exception Program Counter (EPC) register.

(iii) Then the PC jumps to the function which handles
     Overflow.. This function contains the set of
     instructions to overcome overflow.

(iv) Lastly, PC retrives EPC address using mfc0
     (move from coprocessor reg) instruction and
     resume instruction.

" যেই instruction এ overflow হবে PC সেই instruction এর address hold করবে, Then EPC register এ দিয়ে PC সেই instruction এর address store করবে, Then যেই function Overflow handle করবে সাধারণ করে direct সে function এ jump করবে, Then সে function এ কাজ শেষ হলে mfc0 instruction use করে EPC থেকে সেই instruction এর address
Overflow ঘটিয়েছিল তার retrive করে PC তে দিবে, Then বাকি কাজ normally resume হতে পারবে। "


# EPC তে যেই instruction এর address থাকে সেই instruction এ overflow occur হয়/হয়েছিল।


# Coprocessor থেকে কিছু নেওয়ার জন্য mfc0 use করা হয়

## Multiplication:

2 টা approach.

- Long Multiplication Approach
- Fast Multiplication Approach

## Long Multiplication Approach:

- Length of product is the sum of operand lengths

- 32 bit Arch এর জন্য,

  { - Multiplicand and Result এর জন্য 64 bit

  { - Multiplier 32 bit এর জন্য

  → so An, 64 bit একসাথে রাখা যায় না register এ,

    so 2 টা register লাগবে, একটা register এ

    MSB রাখা; আরেকটায় LSB রাখা.
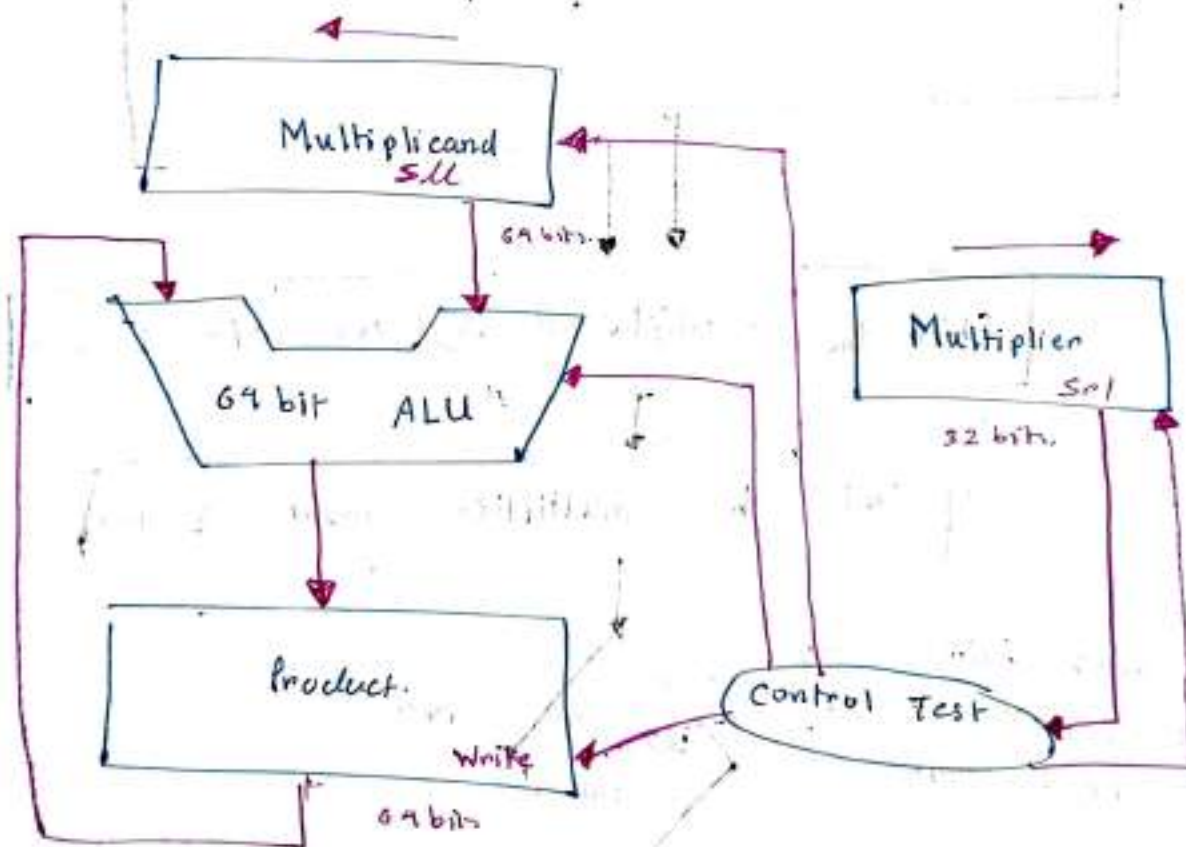
# Multiplicand এর size হবে arch এর size এর double

Example: [Considering 4 bit arch.]

```
     1000  ────▶ Multiplicand
  X  1001  ────▶ Multiplier
  ──────────
     1000
    0000
   0000
  1000
  ──────────
  1001000  ═══▶ Product
```

আমরা product এর
যেই এক register তাতে
এ maximum 8 bit hold
করতে পারে,
এখানে product এর
size 7 bit so, এই
অংশটি সমস্যা হবেনা

# Control Test @ ALU তে Control Signal পাঠায়

Operation করার জন্য,



# Product initially 0 থাকে,

# Number of iteration = number of bit in multiplier =
যত bit এই arch. follow করছে তত bit.

Start

Test Multiplier LSB

LSB = 1

LSB = 0

Add multiplicand to product and place the result in Product register

Shift the multiplicand register left 1 bit

Shift the multiplier right 1 bit

32nd repetation

No

Yes

Done

32 bit anch follow করছে।
so 32nd repetation
form check করছি।

অতএব, 4 bit anch
হলে আমরা আগে check
করবো 4th repetation
কিনা।

Example: [ 4 bit arch.]

▨  8 × 9

⇒ Multiplicand = 8 = .$\underline{0000\ 1000}$,

as 4 bit arch. follow करता.
So, multiplicand length is the double of arch. bit.

Multiplier = 9 = 1001.

| Iteration | Multiplicand 0000 1000 | Multiplier 1001 | Product 0000 0000 |
|---|---|---|---|
| 1 | 0000 1000 | 1001 | 0000 1000 |
|  | 0001 0000 | 100 1 | 0000 1000 |
|  | 0001 0000 | 0100 | 0000 1000 |
| 2 | 0001 0000 | 0100 | 0000 1000 |
|  | 0010 0000 | 0010 | 0000 1000 |
| 3 | 0100 0000 | 0010 | 0000 1000 |
|  | 0100 0000 | 0001 | 0000 1000 |
| 4 | 0100 0000 | 0001 | 0100 1000 |
|  | 1000 0000 | 0001 | 0100 1000 |
|  | 1000 0000 | 0000 | 0100 1000 |

∴ $(0100\ 1000)_2 = (72)_{10} = (8) * 9$

# এর যদি $8 * (-9) = -72$ হয় তাহলে সাইনটা করলাম

Then 8 এর ও এর positive value নিয়ে কাজ

করলাম। Last এ যে final result অাসবে তার

2's complement করে দিবো, Simple.

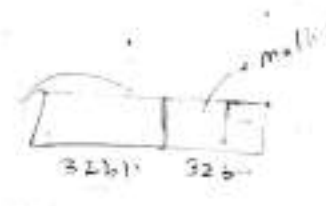## Optimized Multiplier

→ Performs Steps in parallel: add / shift.

⇒ One Cycle per partial-product addition

- That's ok if frequency of multiplication
  is low.

⇒ যেমন, multiplicand ও multiplier 32 bit
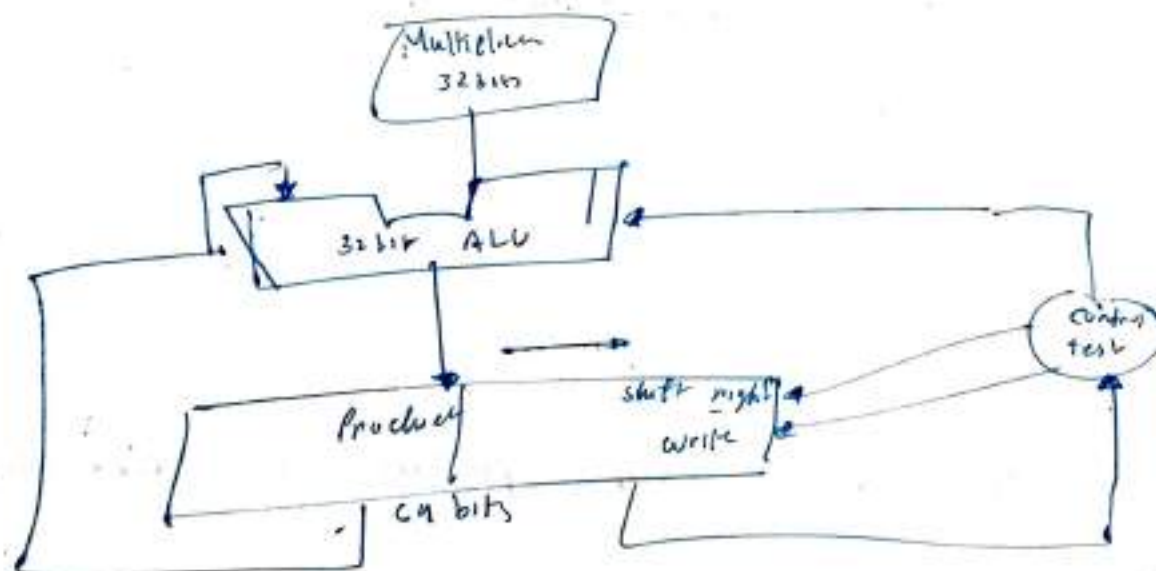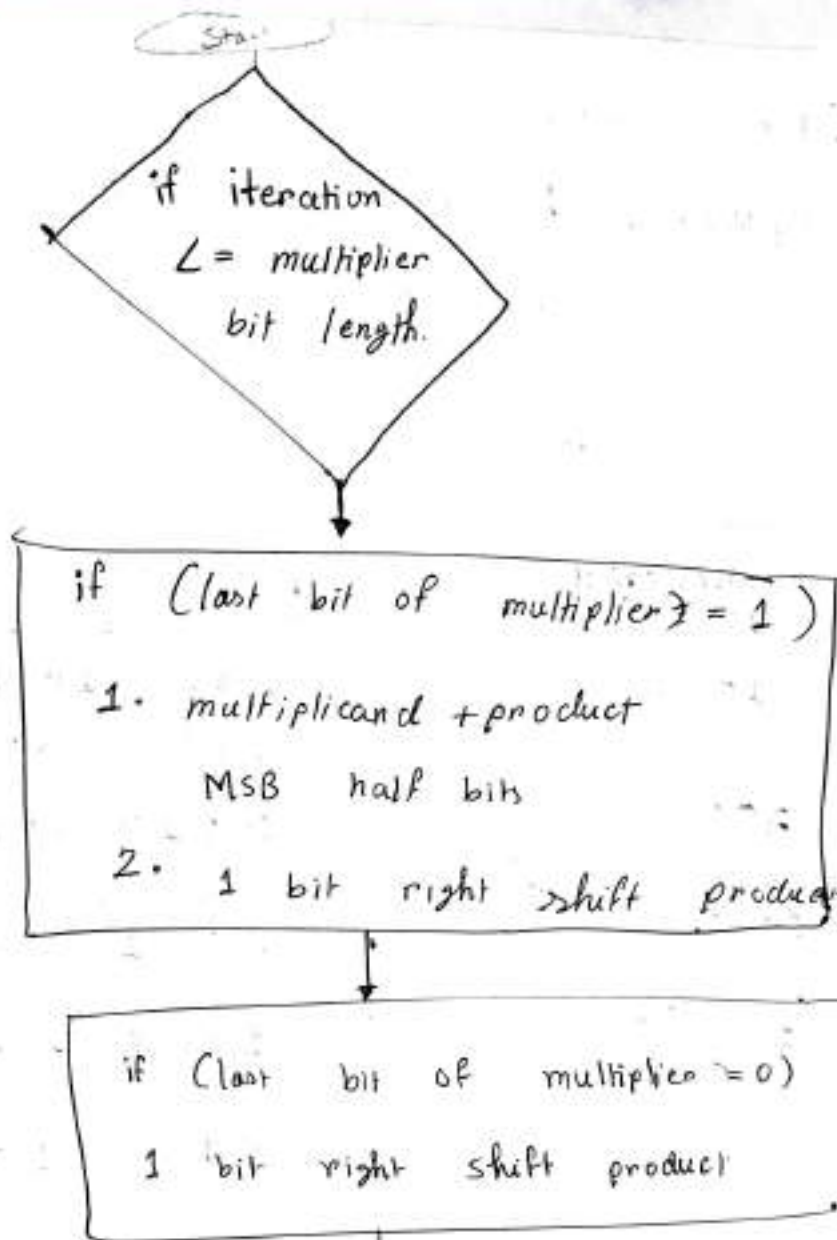
just product 64 bit.

⇒ Product তা 64 bit রেজিস্টার নিবো,



- LSB 32 is used for multiplier

- MSB 32 bit এ Multiplicand এর
  সাথে ALU তে যোগ যে কাজ হয় তার
  উত্তর

**Flowchart:**

Start

if iteration
$\angle$ = multiplier
bit length.

if (last bit of multiplier} = 1 )

1. multiplicand + product
   MSB half bits

2. 1 bit right shift product

if (last bit of multiplier = 0)

1 bit right shift product

Multiplier
32 bits

32 bit ALU

Product

Control test

shift right

write

cn bits

## ▨ Example: [4 bit arch]

☐ Perform multiplication between 8 and 9 using the optimized multiplication approach?

⇒

| Iteration | Multiplicand 1000 | Product |
|---|---|---|
|  |  | 0000    1001 |
| 1 | 1000 | 1000    1001 |
|  |  | 0100    0100 |
| 2 | 1000 | 0010    0010 |
| 3 | 1000. |  |
|  |  | 0001    0001 |
| 4 | 1000 | 1001    0001 |
|  |  | 0100    1000 |

# Throughout the iteration multiplicand া রকোনো

change হয় না.

# efficient কোনো way use করছে সূচীর
তোথি বলৈ করো optimized.

# 32 bit arch এ আমি product এ size 64
bit.

☑ Perform Multiplication between .0110 (Multiplicand) and .110 (Multiplier) using optimized multiplication approach.

| Iteration | Multiplicand (0110) | Product (0000 0110) |
|-----------|---------------------|---------------------|
| 1 | 0110 | 0000 0011 |
| 2 | 0110 | 0110 0011 |
| | | 0011 0001 |
| 3 | 0110 | 1001 0001 |
| | | 0100 1000 |
| 4 | 0110 | 0010 0100 |

# Faster Multiplication approach multiplication গুলো parallely handle করতে পারে।

# MIPS Multiplication in MIPS Instruction :

64 bit register এ রাখা যায় না, তাই, দুটো register

combine করে রাখি, ২টা register আছে,

- high register      Product reg. এ ২টা reg
- low register      আছে,

HI → High Register → MSB 32 bit hold করে,

LO → Low Register → LSB 32 bit hold করে,

## Instruction :

- mult   rs, rt    বা   multu rs, rt

  64 bit    product    in HI | LO
  
      - ২টা   32 bit reg   multiply করি,

  mult → multiplication signea

  multu → multiplication unsigned.

- mfhi   rd

  ↳ mfhi = move from high register.

  ↳ product এ MSB 32 bit টা rd তে store
      করে,

- mflo   rd.

  ↳ move from low register = mflo

  ↳ product এ LSB 32 bit rd তে store করি.

$$\boxed{\text{mul } rd, \; rs, \; rt}$$

$\rightarrow$ यदि exist করে না-

$\rightarrow$ but যদি সে করলে ISB 32 bit পেল T$\rightarrow$

Stop এ হবে।

## Floating Point Representation

• Normalized form

- Conversion from decimal to floating point and vice - versa.

- Floating point arithmetic operation.

## Floating Point

• $-2.34 \times 10^{-56} \longrightarrow$ Floating Point in Normalized form

• $+0.002 \times 10^{-4}$

• $+987.02 \times 10^9$

} Floating point but not normalized

In binary

$1.xxxxxx \times 2^{yy}$

# Normalized form a decimal point এর বামে প্রথম non zero element থাকলে হবে

; binary number represen

# ভালো করে এ ভালো নিয়ম IEEE 759 format we হয়,

# Normalization.

- Numbers with only one non-zero number befr. decimal point.

→ $5.64 \times 10^{33}$ → Normalize

→ $109.64 \times 10^{33}$ → Not normalizen

↳ Turning it into it's corresponding normalizer form

decimal point কে 2 bit left shift করলে

$$1.0964 \times 10^{33+2} = 1.0969 \times 10^{35}$$

যত bit left shift করবে base এ power এ তত add করে দিবে,

→ $0.00578 \times 10^{3}$ → Not normalizen

↳ Transforming to Normalized,

decimal point এ 3 bit right shift করতে হবে,

So, $5.78 \times 10^{0}$

যত bit right shift করবে base এ power থেকে তত minus.

right এর left এ গেলে left shift

left এর right এ গেলে right shift

Single precision → 32 bit;

double precision → 64 bit

## Decimal to Floating Point Representation.

sequentially follow করতে হবে,

- Convert Decimal to Binary

- Convert Binary to Normalized Binary

- Find out biased exponent

- Find out sign bit and fraction

## IEEE Floating Point Representation:

IEEE floating point এ কিভাবে number কে floating point এ representation করতে হবে,

2 বিভাবে floating point representation আছে।

- 32 bit এ এটা
  └ কে Representation Single Precision

- 64 bit এ এটা
  └ কে Representation Double precision

# 64 bit ও আমরা কত range এর number represent করতে পারি.

- কোনাক কেটে এ জনাক বড় number 64 bit এ represent করার মতো সংখ্যা 32 bit দারা possible হয় না.

- bit length যত বেশি হবে তত ভড় range of number cover করতে পারব.

| S | Exponent | Fraction |
|---|----------|----------|

- S হল Sign bit
  - 0 for positive number
  - (-1) for negative number

### Exponent:

| Single Precision | Double Precision |
|------------------|------------------|
| ⟹ 8 bit (exponent) | ⟹ 11 bit (exponent) |
| ⟹ ∅ 23 bits (Fraction) | ⟹ 52 bits (Fraction) |

# Floating point এ exponent সবসময় biased হয়

# biased exponent unsigned

$$X = (-1)^s \times (1 + Fraction) \times 2^{(exponent - Bias)}$$

exponent:
    ↳ actual exponent + Bias.
- Single Precision Bias = 127
- Double precision Bias = ~~1093~~ 1023

Single Precision (32 bit)

→ Biased exponent

| Sign bit | Exponent | Fraction / Mantsc. |
|----------|----------|--------------------|
| 1 bit | 8 bit | 23 bits |

- Sign bit   O হলে positive number; 1 হলে negative number.

- Exponent:
  - Base এর এর power এই exponent
  - এই exponent টিকে Unsigned binary number হিসাবে 8 bit ধরবো.

  $1.11101 \times 2^{(35)}$
  $1.11101 \times 2^{(8)}$

  8 bit unsigned binary range = 0 to $2^8 - 1$ = 0 to 255

  কিন্তু একটা এর exponent এ 0000 0000 and 1111 1111 is reserved...

  | ∴ Range for biased exponent = 1 to 254 |

  ↳ Exponent এ 1 থেকে 254 পর্যন্ত এর 8 bit bin representation এ exponent এর 8 bit এর range.

✓ **if** biased exponent = n bits then.

$$bias = 2^{(n-1)} - 1$$

∴ Bias =

2

∴ For 8 bit biased exponent $= 2^{(8-1)} - 1 = 2^7 - 1 = 127$

✓ Biased exponent = actual exponent + bias.

यदि सारा कोई भी number की 1 to 259

range से लिए जाने आठवा होर 8 bit

representation करता,

# यदि bias का value 1 बताता रहता को 259

आपका बड़ा 24 आठवा को single precision

bitlength

द्वारा न्यूड रहता जैसा आपका रहता

आपता,

∑ से lowest power -126 द्वारा आपता,

as $-126 + 127 = 1$

तो highest power = 127 द्वारा आपता,

as $127 + 127 = 254$

∴ range for exponent: $2^{-126}$ to $2^{127}$

biased exp ऊ range 2[?] 2ऊ 1 to 254

☑ Convert 50.6749 to 32 bits IEEE - 754 Floating Point Representation

∴ part ऊ ग़ेे
part minimum 10
ऊे multiply ग़ेे

→ 50.6749

Binary of 50 = 110010

Binary of 0.6749 = 1010110011

Binary of
.6749 × 2 = 1.3498 → 1 ↓
.3498 × 2 = 0.6996 → 0 ↓ minimum 10 ग़े
.6996 × 2 = 1.3992 → 1

∴ Binary of 50.6749 = 110010 . 1010 11 00 11

∴ Normalized Binary value = 1.$\dot{1}$00$\dot{1}$0 $\dot{1}$0$\dot{1}$0 $\dot{1}$1 00 $\dot{1}$1 × $2^5$
$\underbrace{\qquad\qquad\qquad\qquad}_{\text{fraction / mantissa}}$

Here, actual exponent = 5.

∴ Bias = $2^{(8-1)} - 1$

∴ exponent bias = 5 + 127 = 132 = 10000100

∴ As positive number so sign bit, = 0.

| 0 | 1000 0100 | 0000 0000 10010 1010 1100 11 |
|---|-----------|-------------------------------|

IEEE - 754 Floating Point Rep of 50.6749

= 0 1000 0100 0000 0000 10010 1010 1100 11

= 0X 4219AB300

# Double Precision (64 bit)

```
        Biased
| Sign bit | Exponent |  Fraction / Mantissa        |
|← 1 bit →|←  11 bit →|←        52 bit           →|
```

As fraction और exponent पर तो पर यहां bit इसके चाहे को जानने range of number चाहें पाये,

→ as fraction, 52 bit ऐ so कैसे bit number represent किया पाती है (या/2)

Exponent 11 bit so represent किया पाती है. number of

- Sign bit 0 का positive 1 का negative

- Exponent

 • 11 bit unsigned binary range = 0 to $2^{11} - 1 = 2047$

 ⇒ exponents 0000 0000 000 एवं 1111 1111 111 is reserved.

 ∴ range for biased exponent is $\leq$ 1 to 2046

 ∴ Bias = $2^{(11-1)} - 1 = 2^{10} - 1 = 1023$

For double precision:

- Biased exponent = Actual exponent + Bias (1023)

∴ actual exponent का minimum value $-1022$ और maximum value 1023.

**Convert $-0.232$ to 12 bit IEEE-754 bit Floating Point Representation, where biased exponent is 4 bits**

→

| sign bit | exponent | Fraction |
|---|---|---|
| 1 bit | 4 bit | 7 bits |

0010 0100

Binary of $0.232 = 0.0011101\,1$

Normalized binary of $0.232 = 1.11011 \times 2^{-3}$

exponent:

actual exponent $= -3$

∴ Bias for 4 bits $= 2^{(4-1)} - 1 = 7$

∴ Biased exponent $= -3 + 7 = 4 = 0100$

Sign bit $= 1$

Fraction $= 1101100$

| 1 | 0100 | 1101100 |
|---|---|---|

$= 0 \times A6c.$

Floating Point (single point) to Decimal Conversion

■ 0x F24 00120

→ Hex to Binary.

1111   0010   0,100   0000   0000   ~~0010~~ 0001   0010   0000

Binary according to format.

1   11100100 0   100 0000 0000 0001 0010 0000

Find out exponent and fraction.

Biased exponent = 111 00 100
Biased exponent in Decimal = 228

∴ Actual exponent = 228 − 127

$$= 101$$

Fraction = 0.100 0000 0000 0001 0010 0000

$$= 2^{-1} + 2^{-15} + 2^{-18}$$

$$= 0.5000 343323$$

✓ Decimal value = $(-1)^{\text{sign bit}} \times (1 + \text{Fraction}) \times 2^{\text{exponent}}$

$$= (-1)^1 \times (1 + 0.5000 343323) \times 2^{101}$$

$$= ~~-66~~ \quad 3.8030 388 43 \times 10^{30}$$

∴ Upto 5 decimal point with rounding

$$= 3.80304 \times 10^{30}$$

# Floating Point Addition / Substraction

35. 23 1 42 + 0.000 53

→

X = 35.23142                                    Y = 0.000 53

$X_{bin}$ = 100011 . 00 1110 1111            $Y_{bin}$ = 0.000000000 0.1000 1011

X (Binary Normalized)                          Y (Binary Normalized)

= 1. 00011 00 111 0 1111 $\times 2^5$        = 1. 000 1011 $\times 2^{-11}$

এখানে 5 bit left shift করা
করতে পারি;

∴ Now updated y in Binary normalized
form

= 0.0000 00000 00000000 1000 10 11 $\times 2^5$

**Rule:**
Match the lower exponent
with higher exponent

যাদের মান সমান করতে
তাদের মান সমান power
সমান... তিনি আকারে
তার মান আলাদাভাবে
করতে হবে.

∴ X+Y = ( 1.000  11 00 III 0 1111 $\times 2^5$) + (0. 0000 0000 00 00000 1000 1011 $\times 2^5$)

= (1.000 11 00 III 0 1111 + 0.0000 0000 00 00 000 1000 1011) $\times 2^5$

= 1.000 11 00 III 0 1111 000 1011  $\times 2^5$

= 100011. 00 III 0 111 11 000  1011

= 35. 23 9 222 9121  (Decimal)

## Floating Point Multiplication

☑  $5.234 \times (-0.003)$

⇒  $x = 5.234$  $\qquad$ $\dot{y} = 0.003$

X in Binary = $10L.00\ 111\ 0\ 1111$ ;  Y in Binary

$\quad \xi_o = 0.0000000\ 11000\ 100$
$\quad 101$

$\quad \xleftarrow{} = 1.1000\cdot100\ 101 \times 2^{-9}$.
$\qquad$ (normalized)

x  (Binary normalized)

$= 1.01\ 00\ 111\ 0\ 1111 \times 2^{2}$

# Multiplication को ठारण base को power same कर्ठा।
आगार ता।
• as ऊपर/2 ans negative आएगा ठ12

$X*Y = -(1.01\ 00\ 111\ 0\ 1111 \times 2^{2}) \times (1.1000\ 1\ 00\ 101 \times 2^{-9})$

$= -(1.01\ 00\ 111\ 0\ 1111 \times 1.1000\ 100\ 101) \times 2^{2-9}$

$= -(1.0100\ 111\ 0\ 1111 \times 1.1000\ 100\ 101) \times 2^{-7}$
$\qquad +$

$= -10.000000\ 101 \times 2^{-7}$

$= -0.00000\ 10\ 000000\ 101 \times 2^{0}$

$= -0.01570\ 12939$  (Decimal)

## Floating Point Arithmetic

🔲  51 5 00000  —  BA 10 A 000

→ X = 515 00000

X (Binary)

| Sign bit 1 bit | exponent 8 bit | Mantissa 23 bit |

= 0101  0001  0101  0000 0000 0000  0000  0000

Find out fraction and exponent:

Biased exponent = 101  00010  = $(162)_{10}$

∴ For exponent = ~~127    168~~    162 - 127  = 3.5

∴ Fraction or Mantissa = 0.101  0000  0000  0000  0000  0000

∴ X in Binary Normalized = 1. Fraction × $2^{exponent}$

∴ X in Binary Normalized = 1.101 0000 0000 0000 0000 0000 × $2^{+35}$

Similarly,

Y = BA10A 000

Y (Binary) = 1011  1010  0001  0000  1010  0000  0000  0000

∴ Biased exponent = 011 10100  = 116

∴ exponent = 116 - 127 = -11

∴ Fraction = 0.001  0000  1010  0000  0000  0000

∴ Y in Binary normalized

= -1.001 0000 1010 0000 0000 0000 × $2^{-11}$

[As sign bit = -1]

Again,

X (Binary Normalized)

= 1.101 0000 0000 0000 0000 0000 × $2^{35}$

Y (Binary Normalized)

= − 1.001 0000 1010 0000 0000 0000 × $2^{-11}$

= − 0. [45 0s] 1001 0000 1010 0000 0000 0000 × $2^{35}$

X − (−Y) = X + Y

= (1.101 0000 0000 0000 0000 0000 +

0. [45 0's...] 1001 0000 1010 0000 0000 0000) × $2^{35}$

= 1.101 [42 0's] 1001 0000 1010 0000 0000 0000 × $2^{35}$

■ 7AC D 0000 + 5BCA 0000

⇒ X = 7AC D 0000

X (Binary)

= 0111    1010    1100    1101     0000   0000   0000   0000

= 0 111   1010 1   100 1101     0000   0000   0000   0000

∴ Biased   exponent = 1111 0101   = 245.

∴ Exponent = 245 - 127   = 118    For exponent being 8 bit,

$$Bias = 2^{(8-1)} - 1 = 127$$

∴ Fraction / Mantissa = 0.100 1101   0000 0000   0000 0000

∴ X (Binary normalized)

= 1.100 1101   0000 0000 0000 0000 ∧ $2^{118}$

Y = 5BCA 0000

Y (Binary) = 0101   1011   1100   1010   0000   0000   0000   0000

∴ Biased exponent = 1011 0111   ~ 183.

∴ exponent = 183 - 127   = 56

∴ Fraction   or Mantissa = 100   1010   0000   0000   0000 0000

∴ Y (Binary Normalized)

= 1.100 1010   0000 0000   0000   0000   x $2^{56}$

Y in Binary normalized = 0. [62 0s] 100 1010 00000 0000
                                                0000     0000
                                                  $\times 2^{118}$

$\therefore$ X+Y

$= \Bigg( 1.100\,1101\ 0000\ 0000\ 0000\ 0000 + 0.[62\ 0s]\ 100\ 1010\ 0000\ 0000$
                                                        0000\ 0000 $\Bigg)$
                                                                        $\times 2^{118}$

$= 1.1001101\ [54\ 0s]\ 100\ 1010\ 0000\ 0000\ 0000\ 0000\ \times 2^{118}$

☑ Suppose, X = 19.454 and Y = 3.0124, perform X * y

   using IEEE floating-point representation.

→  X = 19.45 4

X (Binary) = 10011. 0 111 0100

X (Normalized) = 1. 0011 0111 0100 $\times 2^7$.

Y = 3.0124

Y (Bin) = 11. 000 000 1100

Y (Binary Normalized) = 1. 1000000 1100 $\times 2^1$

$X * Y = (1.0011\,0111\,0100 \times 1.1000000\,1100) \times 2^5$

$= 1.1101\,0100\,1011\,0010111\,0000 \times 2^5$

$= 111010.0010\,11\,0010111\,0000$

$= 58.58734$

$$1.0011\,0\,111\,0\,100$$
$$\times\ 1.1000\,0001\,100$$
$$\overline{0\ 00000\ 0000\ 000}$$

$$1\ 00\ 110\,1110\,100\,0^{xx}$$
$$1\ 001\ 10111\,0100\,xx\ ^x$$

$$\overline{1.1101\,01001\,011\,001\,0111\ 0000}$$

$(?. T. 0)$

# MIPS Division:

- Use HI/LU register for result.
  - HI: → 32 bit remainder (উপরাংশ)
  - LO: → 32 bit quotient (ভাগফল)

- Instructions.
  - div rs, rt / divu, rs, rt.
  - No overflow on divide by 0 checking
    - Software must perform check if required
  - Use mfhi, mflo to access result.

# div rs, rt এর, rs register এর value divided by rt.

# div $t0, $t1

■ ভাগের register এর value even অথবা odd এ। ভাগের easily div instruction দিয়ে হবে ভাগ মাঝি।

# যদি মাঝি remainder বা ভাগ করতে চায় বা ভাগের register এ নিতে চায় আমরা, mfhi use করবে,
  
  example `mfhi $t0` → remainder এর value $t0 তে নিয়ে আসবি।

যেকোনো সংখ্যা $10 এর value 0 হয় তাহলে even

তাহলে সংখ্যা odd

# Floating Point Registers and Instruction.

**#** Register এ. যে data রাখতে পারি তা সকলেই int হয়,

**#** Floating Point অন্য Register এ রাখতে পারি না.

**#** Floating Point এর জন্য আলাদা একটি register file

আছে,

> **#** Floating Point Register গুলো
> $fo, $f1, ....., ইত্যাদি.

Named ☐ Co - processor 1.

↳ এটা both int এবং float type data hold

করে,



**#** Co processor 1 এ 32 bit register থাকে এবং

32 টি register থাকে.

Coprocessor 1 এর register file হলো:

$$\$fo \quad থেকে \quad \$f31.$$

■ ► A[3] = (float) x

Where x is in $fo and base address of A in in $s1.

⇒ **swc1** **$fo, 12($s1)**

■ Float (a) = B[2];

Where a is in $fo and base address of B in $s1.

=> lwc1$ fo, 8($s1)

► 

# • Co-processor এ register overwrite করার option নাই,

- ## MIPS Instruction For Single Precision:

⇒ Single Precision 32 bit

$$\$f0 = \$f1 + \$f2$$

add.s $\$f0, \$f1, \$f2$

$$\$f0 = \$f1 - \$f2$$

sub.s $\$f0, \$f1, \$f2$

"add.s যাও sub.s indicate করে single Precision

এর মান।

## MIPS Instruction For Double Precision:

⇒ Double Precision 64 bit.

$$\$f0 = \$f1 + \$f2$$

add.d Valid না, as Coprocessor]

এর 32 bit register গুলো তো

তা 32 bit hold করতে পারে।

but double precision এর 64 bit, যেমন : এ example

এ $\$f0, \$f1, \$f2$ : 32 bit hold করতে but

আসলে আমরা 64 bit. So, Source & Destination

$$\$f0 = \$f1 - \$f2$$

co-processor 2

| fo |
|----|
| f1 |

|  |

64 bit.

যদি $\$f0$ তে আমি একটা 64 bit রাখতে so

২টা slot লাগবে. So, fo তে store

হল v.করা আসলে fo তে f1.

So, fo and f1 combination.

করলে 64 bit হবে.

So, Valid example.

$$\$f_0 = \$f_2 + \$f_4$$

add.d  $\$f_0$, $\$f_2$, $\$f_4$

" add.d বা sub.d indicate করে double precision..

$$\$f_0 = \$f_2 - \$f_4$$

sub.s  $\$f_0$, $\$f_2$, $\$f_4$

## Moving Values from registers:

mtc1  $\$S1$,  $\$f_0$  ⟶  Co processor যেখানে right hand side এ থাকে,

mfc1  $\$S1$,  $\$f_0$

Convert a Register Value from float to int
[Single Precision]

cvt.w.s.  $\$f_0$,  $\$f_1$

Convert করে Coprocessor R' এর মধ্যে,

float value টা r.h.s এ, so, r.h.s. থেকে convert করে l.h.s এ assign.

Convert a register's value from int to float
[Single precision]

left a converted value

cvt.s.w  $\$f_0$,  $\$f_1$

# Double Precision to Single Precision

cvt . s.d.    $ f0,    $ f1

# Single Precision to Double Precision

cvt . d.s.    $ f0,    $ f1.

📝 Float (x) = (int) y + float (z)

int (y) = (float) x - (float) z

Suppose x, y, z are in $f1, $s1 and $f2 correspondingly.

⇒ **MIPS Code:**

mtc1  $s1, $f0  # reg अब mem ओ (int) y आया,

cvt.s.w.  $f0,  $f0  # int y इसका float (y), ও convert

add.s $ f1,  $f0,  $f2  # float(x) = float (y) + float(z)

sub.s $ f3 ,  $f1,  $f2  # float (y) = float (x) - float (z)

cvt.w.s  $f3, $f3  # float(y) इसका int (y) ও convent

mfc1  $s1, $f3  # int (y) = float (x) - float(z)

and int(y) is value is again pushed in the saved register (register file)

- n bit system's unsigned number range

"
$$-2^{(n-1)} \quad \text{to} \quad +2^{(n-1)} - 1$$
"

- n bit systems signed number range

$$0 \quad \text{to} \quad 2^n - 1$$

- Overflow: [ কোন bit arch. follow করবে তার উপর নির্ভর করে।]

| Int: Addition | Int Substraction |
|---|---|
| - Overflow if out of range. | - Overflow if out of range. |
| - If sign bit is 1 while adding two positive number. | - If sign bit is 0 while substracting +ve number, from -he number. |
| - If sign bit is 0 while adding two negative number | - If result sign is 1 while substracting a -ve number from a positive number. |
| - No overflow while adding a positive number and a negative number. | - No overflow while substracting two positive number. or while substracting two negative number. |

# Sign bit = 0 মানে positive

# Sign bit = 1 মানে negative.

- Long Multiplication
  - যদি n bit arch. হয় then,
    - multiplier n bit
    - multiplicand তবে product 2n bit
  - Flow chart for procedure. ⭐

- Optimized Multiplication
  - যদি n bit arch হয় তাহলে;
    - multiplier and multiplicand n bit
    - producta 2n bit.
  - But there's a catch,
    for example in 32bit arch; optimized
    multiplication 1.
      - multiplier 32 bit, multiplicand 32 bit
      - Product 64 bit
        ↳ product এ 64 bit 2 ভাগে ভাগ
        HI ← MSB 32 BIt ——→ denotes ⌣
        LO ← ─ LSB 32 BIt ——→ Multiplier
                              ↳ MSB 32 bit ও multiplicand
                              এ আসে ALV এ ঢুকে
                              কাজ হয় আর হয়

- MIPs Multiplication Instruction.

  → mult rs, rt ; multu rs, rt.

  ⇒ mfhi rd:
    ↳ Product এর MSB 32 bit এটা rd তে store করে.

  ⇒ mflo rd:
    ↳ Product এর LSB 32 bit এটা rd তে store করে.

- Decimal Point
  - left shift করলে power বাড়ে-
  - right shift " " " কমে.

- Conversion From Decimal to Floating Point.

  ① Convert Decimal to Binary

  ② Convent Binary to Normalized Binary.
    → Normalized form: $1.xxxxxx \times 2^{yy}$

  ③ Find out the Biased exponent.

  ④ Find out sign bit and fraction.

  $$X = (-1)^{\text{sign bit}} \times (1 + \text{Fraction}) \times 2^{(\text{Biased exponent} - \text{Bias})}$$

# Floating Point Format.

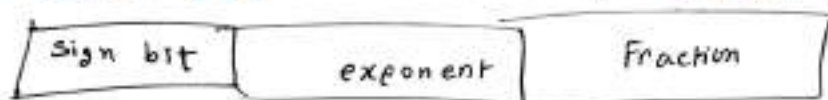- Single Point Precisiun

  ⇒ Sign bit = 1 bit

  ⇒ exponent = 8 bit
     (biased)

  ⇒ Fraction / Mantissa = 23 bit

  ∴ Total    32 bits
     ⇒ Bias = 127

- Double          precision

  ⇒ Sign bit = 1 bit

  ⇒ exponent = 11 bit
     (biased)

  ⇒ Fraction / Mantissa = 52 bih

  ∴ Total  64 Bit
     ⇒ Bias = 1023

| Sign bit | exponent | Fraction |
|----------|----------|----------|

## Formula:

- If exponent field length = ~~8 bit~~ n bit, then,

  $$Bias = 2^{(n-1)} - 1$$

- Biased exponent = Actual exponent + Bias

\# Single Precisiun. exponent Range = -127 to 127.

\# Double Precision     exponent Range = -1022 to 1023

( রকিবুল্লাহ লিখে নি; তাড়াহুড়ো

প্রেম আমার জ্বালা করা )