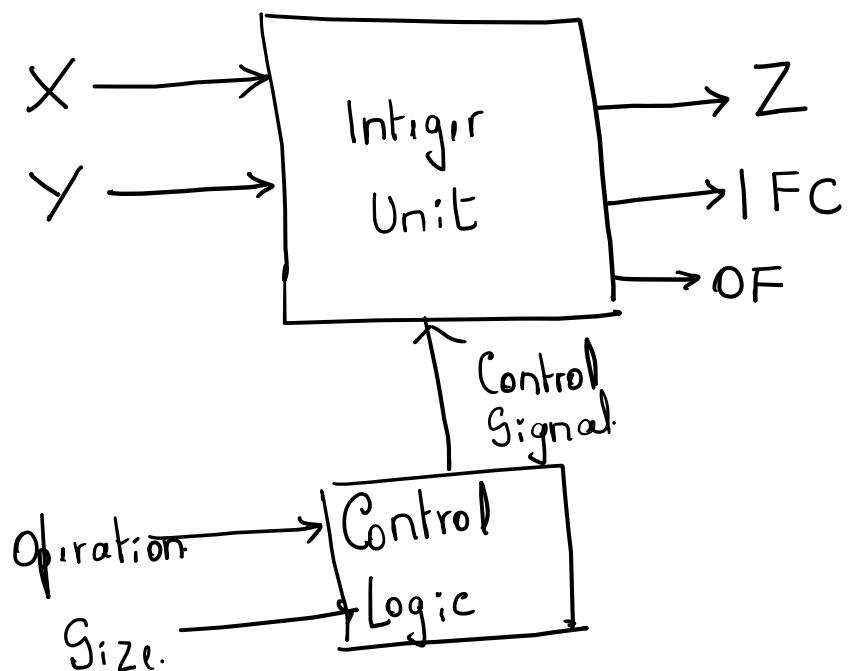


INTEGER UNIT



- Adder / Subtractor
- Unsigned multiplier
- Signed multiplier
- Unsigned divisor
- Signed divisor
- Real valued data

Floating Point Numbers and Operations

Binary format

$$\begin{array}{r}
 2 | 13 \\
 2 | 6 - 1 \\
 2 | 3 - 0 \\
 1 - 1
 \end{array}$$

$$\begin{array}{r}
 0.75 \times 2 \\
 \hline
 1.50 - 1 \\
 \hline
 0.50 \times 2 \\
 \hline
 1.00 - 1
 \end{array}$$

$$13_{(10)} \rightarrow 1101_2$$

$$0.75_{(10)} \rightarrow 0.11_2$$

Fixed Point Representation

- Assume there exists a binary point

Unsigned

$$13 \rightarrow 1101.$$

↑ Implicit binary point to the right of LSB.

$$0.75 \rightarrow \underline{\underline{0.11}}$$

Signed fraction

Positive. $0.75 \rightarrow 0.11$

Negative $-0.\overline{75} \xrightarrow[0]{2}$ 2's complement of 0.11

$$\begin{array}{r} 1.00 \\ \downarrow \\ \hline 1.01 \end{array}$$

$$\underline{-0.\overline{75}_{(10)} \rightarrow 1.01_{(2)}}$$

Generalizing,

$$B = b_0 \cdot b_{-1} \cdot b_{-2} \cdots b_{-(n-1)}$$

Value of B is given by,

$$F(B) = -b_0 \times \overset{0}{2} + b_{-1} \times \overset{-1}{2} + \cdots + b_{-(n-1)} \overset{-(n-1)}{2}$$

Ex:

$$0.\overline{75} \rightarrow 0.11 \quad F(B) = -0 \times \overset{0}{2} + 1 \times \overset{-1}{2} + 1 \times \overset{-2}{2}$$

$$= 0 + \frac{1}{2} + \frac{1}{4}$$

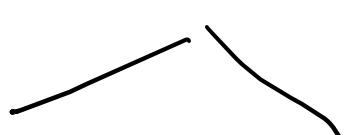
$$= 0.5 + 0.25$$

$$= \underline{0.75}$$

$$-0.\overline{75}_{(10)} \rightarrow 1.01_{(2)}$$

$$\begin{aligned}
 F(B) &= -1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} \\
 &= -1 + 0 + \frac{1}{4} \\
 &= -1 + 0.25 \\
 &= \underline{-0.75}
 \end{aligned}$$

In general \rightarrow n-bits


 i bits can be used for whole part f bits can be used for fraction part.

$$n = i + f$$

$$13.75_{(10)} \rightarrow \underbrace{1101}_{i=4} \cdot \underbrace{11}_{f=2}$$

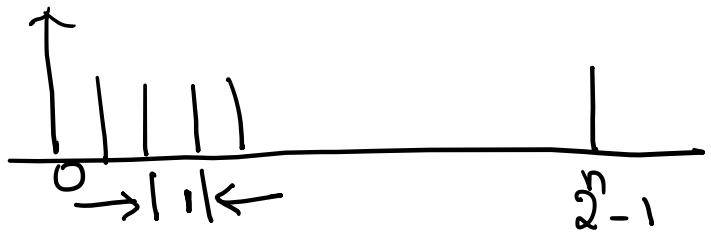
32-bit unsigned fixed-point format

$$n = 32$$

$$l = 32 \quad f = 0 \quad \left\{ \text{Binary point to the right of LSB} \right.$$

Range : 0 to $2^n - 1$

Resolution : 1 $\left\{ \text{Difference between successive numbers} \right\}$



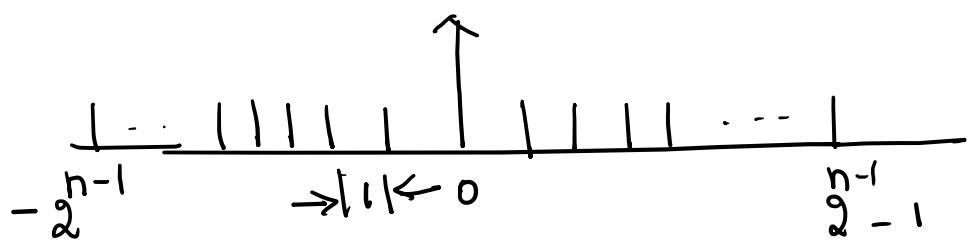
32 bit signed fixed-point

$$n = 32$$

① $i = 32 \quad f = 0$

Range: $-(2^{n-1})$ to (2^{n-1})

Resolution: 1



Range: $-2^{\frac{31}{2}}$ to $\frac{31}{2} - 1$

0 to $\pm 2.15 \times 10^9$

} looks like large
number

Decimal →

② $i = 0 \quad f = 32$: fraction.

- Binary point is assumed to be \rightarrow left of MSB

Decimal nos: $\pm 4.55 \times 10^{-10}$ to ± 1
Very small numbers

$6.0247 \times 10^{23} \text{ mole}^{-1}$ \rightarrow Avogadro's number.

$6.6254 \times 10^{-34} \text{ erg-s}$ \rightarrow Planck's constant

Range \rightarrow 32 bit signed fixed-point representation



May not be sufficient



Scientific computation.

Computer \rightarrow Position of Binary point \rightarrow Variable



Binary point is said to float

Floating point Numbers



Representation.

Range and Resolution

$$n = \text{bits} \quad n = l + f$$

Unsigned

$$\text{Range : } 0 \text{ to } 2^l - 2^f$$

$$\text{Resolution: } \frac{1}{2^f}$$

$$n = 3 \quad \underline{\text{Binary.}}$$

$$\begin{cases} l = 3 \\ f = 0 \end{cases} \quad \text{Range : } 0 \text{ to } 2^3 - 2^0$$

$$0 \text{ to } 2^3 - 1$$

Decimal.

$$0 \text{ to } 10^3 - 10^0$$

$$0 \text{ to } 10^3 - 1$$

$$0 \text{ to } 999$$

$$\text{Resolution: } \frac{1}{2^0} = 1$$

$$10^0 = 1$$

$$\begin{cases} l = 2 \\ f = 1 \end{cases} \quad \text{Range: } 0 \text{ to } 2^2 - 2^{-1}$$

$$0 \text{ to } 10^2 - 10^{-1}$$

$$0 \text{ to } 100 - 0.1$$

$$0 \text{ to } \underline{99.9}$$

$$10^{-1}$$

$$\underline{\text{Resolution.}} \quad \frac{1}{2^1}$$

$$\begin{cases} l = 1 \\ f = 2 \end{cases} \quad \text{Range : } 0 \text{ to } 2^1 - 2^{-2}$$

$$0 \text{ to } 10^1 - 10^{-2}$$

$$0 \text{ to } 10 - 0.01$$

$$0 \text{ to } \underline{9.99}$$

Resolution: $\frac{1}{2}^2$

$i = 0$
 $f = 3$

Range: 0 to $2 - \frac{1}{2}^3$

10^{-2}

0 to $10^0 - 10^{-3}$

0 to 1 - 0.001

0 to 0.999

10^{-3}

Resolution: $\frac{1}{2}^3$

Scientific notation: Decimal.

253.75×10^0

25.375×10^1

2.5375×10^2

0.25375×10^3

0.025375×10^4

Mantissa.

Significand.



Same number.

Exponent

Sign -25.3×10^{-2} → Exponent
Significand
Mantissa.

Negative exponent

$$13.75 \rightarrow 1101.11 \times 2^0$$

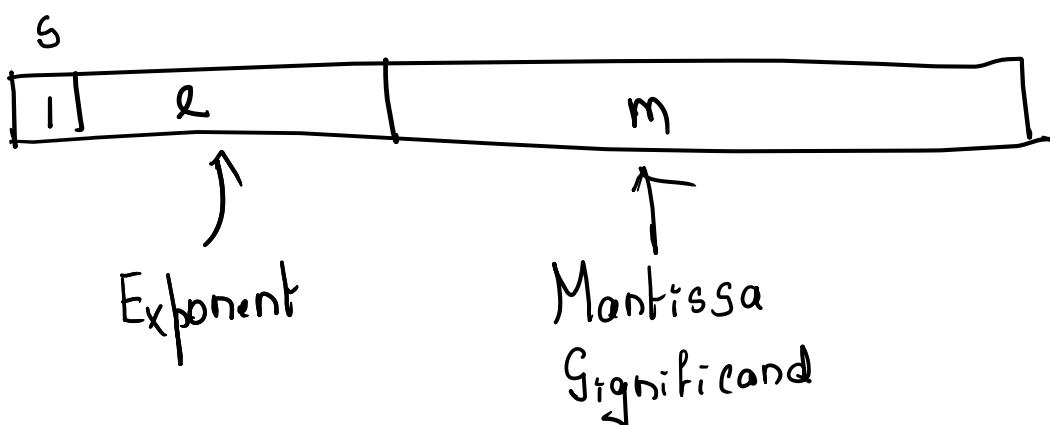
$$110.111 \times 2^1$$

$$11.0111 \times 2^2$$

$$1.10111 \times 2^3$$

$$\underbrace{11011.1}_{\text{Significand}} \times 2^{-1} \quad \text{Exponent}$$

13.75



$$n = s + e + m$$

Normalized representation

- Ensure that there is only one bit to the left of the binary point and that bit is always 1

$$13.75 \rightarrow 1101.11 \times 2^0$$

↓ Normalize.

$$\underbrace{1.10111}_{\text{Significand}} \times 2^{\textcircled{3}}$$

0	$<3>$	10111...
---	-------	----------

$$\begin{array}{r}
 0.875 \times 2 \\
 \hline
 0.750 \times 2 \\
 \hline
 \textcircled{1} \quad 0.500 \times 2 \\
 \hline
 1.000
 \end{array}$$

$$\begin{array}{l}
 0.875 \rightarrow 0.111 \times 2^0 \\
 \downarrow \text{Normalize.} \\
 \underbrace{1.11}_{\text{1}} \times 2^{-1}
 \end{array}$$

0	$<-1>$	1100...
---	--------	---------

Signed number.

Biased exponent

\downarrow
 True exponent in Excess-K format
 1.11×2^{-1} True exponent

0	Biased ex \rightarrow [-1]	Mantissa.
---	------------------------------	-----------