

26-Mar-2024

Term 1 / Week 9

Floating Point Number

Homework Solution

Given $(0.35)_{10} \cong (0.01011)_2$
by ignoring recurring numbers

• Calculate the rounding - off error

$$\begin{aligned}
 & 0.01011 \\
 & \begin{matrix} - & - & - & - & - \\ 2 & 2 & 2 & 2 & 2 \end{matrix} \\
 & = 0 \times 2^1 + 1 \times 2^{-2} + 0 \times 2^{-3} + 1 \times 2^{-4} + 1 \times 2^{-5} \\
 & = 1 \times \frac{1}{4} + 1 \times \frac{1}{16} + 1 \times \frac{1}{32} \\
 & = 0.25 + 0.0625 + 0.03125 \\
 & = 0.34375
 \end{aligned}$$

-3-

Floating Point Number

• Reviewing ^{the} Example from last week, we have

$$(25.625)_{10} = (11001.101)_2$$

• This format is called "Fixed Point" - the fractional value follows the decimal or the binary point.

• For computer storage, it is convenient to 'set' the 'point' at one end. Normally, the 'point' is shifted to the left most end.

$$\begin{aligned}
 \therefore \text{Error} &= 0.35 - 0.34375 \\
 &= 0.00625
 \end{aligned}$$

$$\begin{aligned}
 \% \text{Error} &= \frac{0.00625}{0.35} \times 100 \\
 &= \underline{\underline{1.79\%}}
 \end{aligned}$$

Trying next approximation,

$$(0.35)_{10} = (0.010110011)_2$$

we get

$$(0.010110011)_2 \cong (0.34961)$$

$$\begin{aligned}
 \% \text{Error} &= \frac{0.35 - 0.34961}{0.35} \times 100 \\
 &= \underline{\underline{0.11\%}}
 \end{aligned}$$

-4-

$$\begin{aligned}
 (11001.101)_2 &= 0.11001101 \times 2^{+5} \\
 &= [0.11001101 \times (10)^{+10}]
 \end{aligned}$$

• We have used 2^x for clarity. Note that $(2)_{10} = (10)_2$

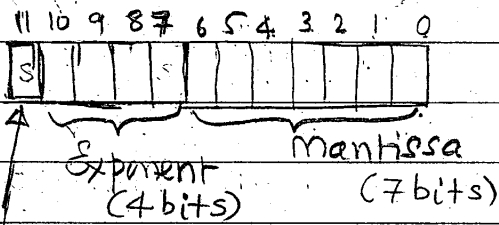
• In the computer, the above value is stored as below:

11001101 is called the Mantissa
+101 is called the Exponent

• Note that ^{the} "Exponent" is always an Integer or a whole number.

• Negative Integers are stored as 2's complement!

Assuming we have a 12-bit machine, the above number is stored as below:



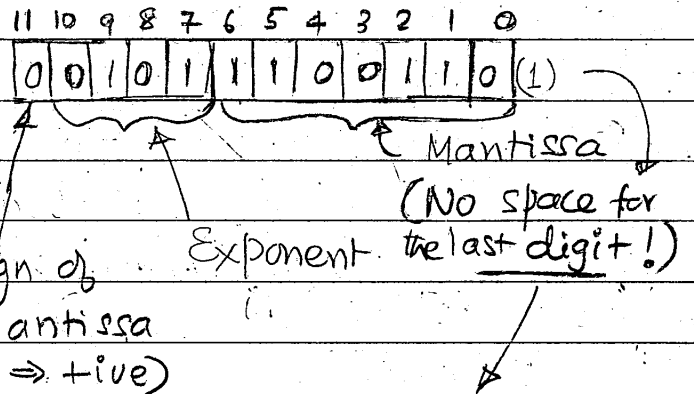
's' indicates the sign bit of the mantissa!

We have: Mantissa \Rightarrow 11001101
Sign of Mantissa \Rightarrow 0 (+ive)

Exponent is stored as a 2's complement for negative values.

We have a positive exponent (+5), hence no need for 2's complement.

Hence, the number is stored as below:



IEEE 754 standard for floating point numbers uses a trick! (Next week!)

Example

Add $25.625 + 4.125$

In floating point form

$25.75 = 0.25625 \times 10^2$
 $4.125 = 0.4125 \times 10^1$

To equate exponents, let us adjust this exponent it should be 10^2 !

$0.4125 \times 10^1 = 0.04125 \times 10^2$

\therefore the sum is

0.25625×10^2
 $+ 0.04125 \times 10^2$
 0.29750×10^2

The answer is 0.29750×10^2
In "fixed point" form the answer is 29.75!

In modern computers, the floating point arithmetic is done by dedicated hardware!

Homework

Add the numbers in the above Example by changing both exponents to 1. $\left. \begin{matrix} \text{which method?} \\ \text{is preferable?} \end{matrix} \right\}$