230201003 Tunahan YADİGARBİGÜN
230201057 Furkan Emre YILMAZ

**Q1 –**

$$\frac{7}{8} = \frac{1}{8} + \frac{1}{4} + \frac{1}{2} = 2^{-3} + 2^{-2} + 2^{-1} \longrightarrow (0.111)_2$$

$$10.5 = 10 + \frac{1}{2} \qquad (10)_{10} = (1010)_2$$

$$(\tfrac{1}{2})_{10} = (0.1)_2 \qquad (10.5)_{10} = (1010.1)_2$$

$$12.8 = 12 + 0.8 \qquad (12)_{10} = (1100)_2$$

$$0.8 \times 2 = 0.6 + 1$$
$$0.6 \times 2 = 0.2 + 1 \qquad (0.8)_{10} = (\overline{1100})_2$$
$$0.2 \times 2 = 0.4 + 0$$
$$0.4 \times 2 = 0.8 + 0 \qquad (12.8)_{10} = (1100.\overline{1100})_2$$

$$\frac{2}{3} \times 2 = \frac{1}{3} + 1$$
$$\qquad\qquad\qquad\qquad (\tfrac{2}{3})_{10} = (0.\overline{10})_2$$
$$\frac{1}{3} \times 2 = \frac{2}{3} + 0$$

$$3.2 = 3 + 0.2 \qquad (3)_{10} = (0011)_2$$
$$0.2 \times 2 = 0.4 + 0 \qquad (0.2)_{10} = (\overline{0011})_2$$
$$0.4 \times 2 = 0.8 + 0$$
$$0.8 \times 2 = 0.6 + 1 \qquad (3.2)_{10} = (0011.\overline{0011})_2$$
$$0.6 \times 2 = 0.2 + 1$$

Find the IEEE doble precision representation $fl(x)$, and find the exact difference $fl(x) - x$ for the numbers $\frac{1}{3}$, 3.3, and $\frac{9}{7}$

a) First convert $\frac{1}{3}$ to binary

$\frac{1}{3} \times 2 = \frac{2}{3} + 0$

$\frac{2}{3} \times 2 = \frac{1}{3} + 1$

$\frac{1}{3} \times 2 = \frac{2}{3} + 0$

$\frac{2}{3} \times 2 = \frac{1}{3} + 1$

$\left(\frac{1}{3}\right)_{10} = (0.\overline{01})_2$ we can check it

$2^2 x = (01.\overline{01})_2$

$1.x = (0.\overline{01})_2$

$\overline{\phantom{3x = (01)_2}}$

$3x = (01)_2$

$x = \frac{1}{3}\wedge$

Second step using IEEE doble precision $fl(x)$,
we first convert number to normalized format

$$(0.\overline{01})_2 = (01.\overline{01}) \times 2^{-2}$$

$$= (1.\overline{01}) \times 2^{-2}$$

we know that in doble precision we have 52 bit for mantissa so, we can write number such as

$$\underbrace{1.0101010101010101 \text{ -------- } 01,}_{\text{52 bit}} \underbrace{0101 \text{ ---}}_{\text{infinite tail}} \times 2^{-2}$$

to represent the number we discard infinite tail

$$2^{-2} \times 2^{-52} \; (0.\overline{01})_2 \rightarrow \text{discarded}$$

$$fl\left(\frac{1}{3}\right) = \frac{1}{3} - \frac{1}{3} \times 2^{-54} \qquad (0.\overline{01})_2 = \frac{1}{3}$$

$$fl\left(\frac{1}{3}\right) - \frac{1}{3} = -\frac{1}{3} \times 2^{-54}$$

Using the same steps as before

b) first convert 3.3 to binary

$$3.3 = 3.0 + 0.3$$

$0.3 \times 2 = 0.6 + 0$
$0.6 \times 2 = 0.2 + 1$
$0.2 \times 2 = 0.4 + 0$
$0.4 \times 2 = 0.8 + 0$
$0.8 \times 2 = 0.6 + 1$

$0.6 \times 2 = 0.2 + 1$
$0.2 \times 2 = 0.4 + 0$
$0.4 \times 2 = 0.8 + 0$
$0.8 \times 2 = 0.6 + 1$

$$3/2 = 1.2 +1 \uparrow$$
$$1/2 = 0.2 + 1$$

$$(3)_{10} = (11)_2$$
$$(0.0\overline{1001})_2 = (0.3)_{10}$$

$$(3.3)_{10} = (11.0\overline{1001})_2$$

Normalized form can be written as

$$(1.10\overline{1001}) \times 2^1$$

can be represented in computer as

$$1.1010011001 \underbrace{\text{------} 100110}_{52 \text{ bit}} \underbrace{\overline{0110011001}\text{---}}_{\text{infinite tail}} \times 2^1$$

to represent the number we discard infinite tail

$$2^1 \times 2^{-52} \times (0.\overline{0110})_2$$

$$fl(3.3) = 3.3 - 2^{-51} \times (0.4)$$

$$fl(3.3) - 3.3 = -2^{-51} \times (0.4)$$

$$(0.\overline{0110})_2 \Rightarrow 2^4 x = (110.\overline{0110})_2$$
$$-\quad x = (0.\overline{0110})_2$$

$$15x = (110)_2$$

$$x = \frac{6}{15} = \frac{2}{5} = 0.4$$

c) First convert $\frac{9}{7}$ to binary

$$\frac{9}{7} = 1 + \frac{2}{7}$$

$$\left(\frac{9}{7}\right)_{10} = (1.0\overline{100})_2$$

$$\frac{2}{7} \times 2 = \frac{4}{7} + 0$$

This number already in normalized form so our multiplier $2^0$

$$\frac{4}{7} \times 2 = \frac{1}{7} + 1$$

$$(1.0\overline{100})_2 \times 2^0$$

$$\frac{1}{7} \times 2 = \frac{2}{7} + 0$$

We can represent in $fl(x)$ as

$$\frac{2}{7} \times 2 = \frac{4}{7} + 0$$

$$\underbrace{1.0100\,100\,\text{-----}\,100}_{52\,bit}\,\overbrace{\overline{100}\,100\,100\,\text{---}}^{\text{infinite tail}} \times 2^0$$

$$\frac{4}{7} \times 2 = \frac{1}{7} + 1$$

we need to round up

$$\left(\frac{2}{7}\right)_{10} = (0.0\overline{100})$$

if we round up our number can be represented as

$$\underbrace{1.0100100\,\text{-----}\,101}_{52\,bit}$$

At the end lets consider what we added what we substracted

we added

$$2^0 . 2^{-52}$$

we substracted or discarded

$$2^0 . 2^{-52} \underbrace{(0.\overline{100})_2} \longrightarrow (0.\overline{100})_2 = \left(\frac{4}{7}\right)_{10}$$

Therefore result is

$$fl\left(\frac{9}{7}\right) \sim \frac{9}{7} + 2^{-52} \times 2^0 - \frac{4}{7} \times 2^{-52} \times 2^0$$

$$fl\left(\frac{9}{7}\right) - \frac{9}{7} = 2^{-52} - \frac{4}{7} \times 2^{-52}$$

Do the following computation by hand in IEEE double precision computer arithmetic using Round the nearest rule: $4.3 - 3.3$

first convert the numbers binary

$$4/2 = 2.2 + 0 \uparrow \qquad (100)_2 = 4$$
$$2/2 = 1.2 + 0 \qquad (0.0\overline{1001})_2 = (0.3)_{10}$$
$$1/2 = 0.2 + 1 \qquad \Bigg\rangle (100.0\overline{1001})_2$$

we already know binary equivalent of $3.3$ is

$$(3.3)_{10} = (11.0\overline{1001})_2$$

Second represent numbers in $fl(x)$

$$(100.01001)_2 = (1.000\overline{1001})_2 \times 2^2$$
$$\underset{\mathit{cf}}{}$$

$$1.000 1001 1001 1001 \text{----} \underbrace{1}_{52 \text{ bit}} \underbrace{\overline{0011001001}}_{\text{infinite tail}} \times 2^2$$

to represent the number we discarded infinite tail

$$= 2^2 . 2^{-52} (0.\overline{0011})_2 \qquad (0.\overline{0011})_2 \Rightarrow 2^4 x = (11.\overline{0011})_2$$
$$\underline{\qquad\qquad x = (0.\overline{0011})_2}$$
$$= 2^{50} \times (0.2) \qquad\qquad 15 x = (11)_2$$

this number represent our error. from question two we know
the error of $3.3$ is $2^{-51} \times 0.4$. $\qquad x = \dfrac{3}{15} = \dfrac{1}{5} = 0.2$

If we need to speak clearly

$$fl(4.3) + 2^{-50} \times (0.2) = 4.3$$
$$fl(3.3) + 2^{-51} \times (0.4) = 3.3$$

If we subtract this numbers because of error we made for each floating print number is same we will find the exactly $1$

$$4.3 - 3.3 = 1$$
$$fl(4.3) - fl(3.3) = 1 \qquad \Bigg\rangle =$$

# Q4 -

Note that $f(x) = x^3 - 9$   $f(2) = -1$   $f(3) = 18$

| $i$ | $a_i$ | $f(a_i)$ | $c_i$ | $f(c_i)$ | $b_i$ | $f(b_i)$ |
|---|---|---|---|---|---|---|
| 0 | 2 | – | 2.5 | + | 3 | + |
| 1 | 2 | – | 2.25 | + | 2.5 | + |
| 2 | 2 | – | 2.125 | + | 2.25 | + |
| 3 | 2 | – | 2.0625 | – | 2.125 | + |
| 4 | 2.0625 | – | 2.09375 | + | 2.125 | + |
| 5 | 2.0625 | – | 2.078125 | – | 2.09375 | + |
|   |   |   | 2.0703125 |   |   |   |
|   |   |   | 2.0859375 |   |   |   |
|   |   |   | 2.08203125 |   |   |   |
|   |   |   | 2.080078125 |   |   |   |
|   |   |   | 2.0810546875 |   |   |   |
|   |   |   | 2.08056640625 |   |   |   |
|   |   |   | 2.080322265625 |   |   |   |
|   |   |   | 2.0802001953125 |   |   |   |
|   |   |   | 2.08013916015625 |   |   |   |
|   |   |   | 2.080106445578125 |   |   |   |
|   |   |   | 2.08009277683790625 |   |   |   |
|   |   |   | 2.08008575435941312 |   |   |   |
| 19 |   |   | 2.080061939396972656 |   |   |   |
| 20 |   |   | 2.08008384704058954 |   |   |   |
|   |   |   | 2.08008289373715820 |   |   |   |
|   |   |   | **2.08008337020087402** |   |   |   |

Since the rule is
obvious, no need
to indicate all the
numbers here

A solution is correct within $p$ decimal places if the error is less than $0.5 \times 10^{-p}$

$$\frac{1}{2^{n+1}} < 0.5 \times 10^{-6}$$

$$n > \frac{6}{\log 2} = \frac{6}{0.301} = 19.9$$   therefore, we applied the method

until $i = 20$.

**Q5 –**   Note that $f(x) = \cos x - \sin x$   $f(\pi) = -1$   $f(0) = +1$

| i | $a_i$ | $f(a_i)$ | $c_i$ | $f(c_i)$ | $b_i$ | $f(b_i)$ |
|---|---|---|---|---|---|---|
| 0 | $\pi$ 1.570796 | – | 1.570796 | + | 0 | + |
| 1 | 1.570796 | – | 0.785398 | – | 0 | + |
| 2 | 1.178097 | – | 1.178097 | – | 0.785398 | + |
| 3 | 0.981747 | – | 0.981747 | – | 0.785398 | + |
| 4 | 0.981747 | – | 0.883572 | + | 0.785398 | – |
|   |   |   | 0.834485 |   |   |   |
|   |   |   | 0.809941 |   |   |   |
|   |   |   | 0.797670 |   |   |   |
|   |   |   | 0.791534 |   |   |   |
|   |   |   | 0.788466 |   |   |   |
|   |   |   | 0.786932 |   |   |   |
|   |   |   | 0.786165 |   |   |   |
|   |   |   | 0.785781 |   |   |   |
|   |   |   | 0.785589 |   |   |   |
|   |   |   | 0.785494 |   |   |   |
|   |   |   | 0.785446 |   |   |   |
|   |   |   | 0.785422 |   |   |   |
|   |   |   | 0.785410 |   |   |   |
| 19 |   |   | 0.785404 |   |   |   |
| 20 |   |   | 0.785401 |   |   |   |
|   |   |   | 0.785399 |   |   |   |

Find each fixed point of $g(x) = x^2 + \frac{1}{2}x - \frac{1}{2}$ and decide whether fixed Point Iteration is locally convergent to it.

§ The real number $r$ is a fixed point of the function $g$ if $g(r) = r$

§§ Fixed-Point Iteration is locally convergent if $|g'(r)| < 1$

Based on this definition lets first find the fixed point of the function

$$x = x^2 + \frac{1}{2}x - \frac{1}{2}$$

This means we are looking for solutions to

$$y = x$$
$$y = x^2 + \frac{1}{2}x - \frac{1}{2}$$

Drawing these functions we can learn whether we have such point or not easily. We will find two intersection point $-0.5$ and $1$ these are our fixed points for this functions.

Lets try to find these points are locally convergent or not

$$g'(x) = 2x + \frac{1}{2}$$

$$g'(-0.5) = -1 + \frac{1}{2}$$

$$= -0.5$$

So $|g'(-0.5)| < 1$ so $-0.5$ point is locally convergent

$$g'(1) = 2 + \frac{1}{2}$$

$$= \frac{3}{2}$$

so $|g'(1)| < 1$ is not satisfied. It means, it is not locally convergent for $1$.

When we apply fixed-point iteration method for different initial point we can see function converges to $-0.5$ sufficiently near this point. For example;

$x_0 = 0.9$      0.9
                0.76
                0.4576
                -0.06180224 ....
                -0.052708160 ---
                ¦
                -0.5

Q8 - Forward error : $|r - x_a| \longrightarrow |0.75 - 0.74| = 0.01$

Backward error : $|f(x_a)| \longrightarrow |f(0.74)| = 0.0016$


Q9 -

$x_0 = $ initial guess

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad \text{for} \quad i = 0, 1$$

$$f(x) = x^3 + x^2 - 1 \qquad f'(x) = 3x^2 + 2x$$

$$x_0 = 1$$

for $i = 0$,

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} = 1 - \frac{1}{5} = 0.8$$

for $i = 1$,

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)} = 0.8 - \frac{0.152}{3.52} \approx 0.7568$$