# Sept -16

## Floating point operations

### Addition

$$2.3$$
$$10.764 \quad +$$

$$
\begin{array}{r}
10.764 \\
+ \ 2.3 \\
\hline
13.064
\end{array}
$$

1) Align the decimal points.

$2.3 \times 10^2 + 4.6$

$= (0.046 \times 10^2)$

$\begin{array}{r} 2.3 \\ + 0.046 \\ \hline 2.346 \end{array}$

1) Align Decimal Points

    a) Shift the smaller number to the right s.t both the exponents match

$\underset{A}{(1.x)} + \underset{B}{\left(\begin{array}{c} 1.y \\ 0.z \end{array}\right)}$ or

2)

    A + B

    (Perform the addition $\Big\}$ [Maintain some extra precision beyond 23 bits]

3)

$$(1.1)_b + (1.1)_b = (11.0)_b$$

Not ↑ Normalized

Renormalize the number.

$$(11)_b = (1.1 \times 2^1)_b$$

4) At the moment, number is

$$\left[ 1.(x_1 \underbrace{- - - - - x_{23}}_{23 \text{ bits}}) \times 2^E \right]$$

$$1. \underbrace{x_1 - - - - x_{23}}_{} \; (\overbrace{\phantom{---}}^{}) \\ (x_{24} - - - -)$$

Discard digits $(x_{24})$ onwards

Try to round $x_{23}$

Example:

$1. \left( \underbrace{1 \quad 1 \quad 1 \quad \cdots 1}_{23} \right) \quad \underbrace{(1 1 1 1)}_{4}$ ✗

$= (10.0)_b$

$= (1.0)_b \times 2^1$

5) After rounding — my number might get un-normalized.

Renormalize.

## Rounding

Assume that you want to round to the next integer.

$$4.5 =$$
$$4.4 = 4$$
$$4.6 = 5$$

Rounding Policies:

Round - Up : $4.1 \rightarrow 5$
(+∞)

Round - Down : $4.9$
(-∞)  $\rightarrow 4$

Truncate : $5.369 \rightarrow 5$

Round
(complicated)

$(+)ve$    Trunc $=$ Round-Down

$(-)ve$    Trunc $=$ Round-up

$1.0 \ldots . 1.5 \rightarrow 1$

$1.50001 \ldots . 1.9999$

$\rightarrow 2$

complicated

### Assume

Implement round to next integer
using the complicated rounding
scheme.

$$N = \left\{ 1 . r \underbrace{x_1 \ldots . x_n}_{n} \right\}_b$$

$(r == 0)$

$N = 1$

$(r == 1) \{$

$\quad S = x_1 / x_2 \ldots x_n$

$$1 \cdot r \underbrace{(x_1 | x_2 \cdots}_{(x_1 \cdots x_n)} | x_n) = s$$

round bit

sticky bit

```
if (s==0) {
    N = 1
}
else if (s==1) {
    N = (2)_d = (10)_b
}
```

Result:

$$1. \; [\underline{x_1} \; \text{----} \; x_{23}] \; [\boxed{x_{24}} \; \underbrace{x_{25} \cdots x_{30}}]$$

$r$

$OR$

$s$

We decide whether to increment $x_{23}$

# Extra

## Addition

✓

## Subtraction

1) Take a look at sign bits.

2) Do a sub or add

3) Set value of final sign bit

## Division

1) Sign of the result

## Multiplication

1) Figure out sign of the result

2) Add the exponents

2) Subtract exponents

# Floating Point Tricks

Is:

$(y > 0)$

$(x + y > x)$

(maybe)

$\left\{ a + b - c \times d + e \right\}$

$$\left. \begin{array}{l} 1.0000 \\ + 1 \times 2^{-50} \end{array} \right\} \text{Use double}$$
$$\text{(52 bits)}$$

$$\overline{1.0000\ldots}$$

$$1 + 2^{-50}$$

FP operations are not associative