# Floating-Point Addition

CS520

Dept. of Computer Science
Univ. of New Hampshire

$$2.34 \times 10^2$$
$$+ \; 2.56 \times 10^0$$
_____

$$2.34\boxed{00} \times 10^2$$
$$0.0256 \times 10^2$$
_____
$$2.36\boxed{56} \times 10^2$$

$$\boxed{2.37 \times 10^2}$$

# Floating - Point Addition

1. shift significand of the smaller number to the right until exponents agree

2. add the significands

3. normalize the sum & check for overflow/underflow

$$\overline{10.0\cdots}$$
$$\rightarrow$$

$$\overline{0.00001\cdots}$$
$$\leftarrow$$

4. round the sum
   $\hookrightarrow$ might then need another normalization

7F 2A A A A A
+ 78 F8 78 78

0.111 1111 0010 1010 1010 1010 1010 1010

0.111 1000 1.111 1000 0111 1000 1111 1000

0000 1101

$\Delta_{exp} = 13_{10}$

1.010 1010 1010 1010 1010 1010
+ 0.000 0000 0000 0111 1100 0011 | 11 | 000 1111 1000
1.010 1010 1011 0010 0110 1101
+1

1.010 1010 1011 0010 0110 1110

Guard bits

sticky bit = 1

0111 1111 0010 1010 1011 0010 0110 1110 ✓

7 F 2 A B 2 6 E ✓

# limitations of floating-point numbers

remember: they are <u>approximate</u>!

so: $x \not= y$ is problematic

$$|x - y| < \varepsilon$$

## also problematic

$$1.2345 \times 10^3$$
$$-1.2341 \times 10^3$$
$$0.0004 \times 10^3$$

$\leftarrow$

$$4.0000 \times 10^{-1}$$
?

really just garbage

and floating-point addition is not associative

$$-1.5 \times 10^{38} + (1.5 \times 10^{38} + 1.0) =$$
$$-1.5 \times 10^{38} + 1.5 \times 10^{38} =$$
$$0.0$$

$$(-1.5 \times 10^{38} + 1.5 \times 10^{38}) + 1.0 =$$
$$0.0 + 1.0 =$$
$$1.0$$

To learn more :

take a course in Numerical Analysis

e.g. Math 753

For one disaster story :

google "Ariane 5"