

# CompOrg HW2:

3.13.  $62 \times 12 \rightarrow 00111110 \times 1100$

Iter	Step	Multiplier	Multiplicand	Product
0	Initial values	1100	0011 1110	0000 0000
1	1a: 0 $\Rightarrow$ no operation	1100	0011 1110	0000 0000
	2: 01 multiplicand	1100	0111 1100	0000 0000
	3: 01 multiplier	0110	0111 1100	0000 0000
2	1a: 0 $\Rightarrow$ no op	0110	0111 1100	0000 0000
	2: "	0110	1111 1000	0000 0000
	3: "	0010	1111 1000	0000 0000
3	1a: 1 $\Rightarrow$ prod += mcard	0011	1111 1000	1111 1000
	2: "	0011	1111 0000	1111 1000
	3: "	0001	1111 0000	1111 1000
4	1a: 1 $\Rightarrow$ prod += mcard	0001	1111 0000	1110 1000
	2: "	0001	1110 0000	1110 1000
	3: "	0000	1110 0000	1110 1000

Final product  $\rightarrow 11101000 = 74_{10}$

3.18.  $74 \div 21 \rightarrow 01001010 \div 010101$  (b-bit)

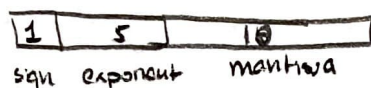
iter	Step	Quotient	Divisor	Remainder
0	initial values	000000	01010100	01001010
1	1: rem = rem - div.	000000	01010100	01110110
	2b: rem < 0 $\Rightarrow$ +Div, shift, Q0=0	000000	01010100	01001010
	3: shift div right	000000	00101010	01001010
2	1: rem = rem - div	000000	00101010	00100000
	2b: rem > 0 $\Rightarrow$ shift, Q0=1	000001	00101010	00100000
	3: shift div right	000001	00010101	00100000
3	1: rem = rem - div	000001	00010101	00001011
	2b: rem > 0 $\Rightarrow$ shift, Q0=1	000011	00010101	00001011
	3: shift div right	000011	00001010	00001011
4	1: rem = rem - div	000011	00001010	00000001
	2b: rem > 0 $\Rightarrow$ shift, Q0=1	000111	00001010	00000001
	3: shift div right	000111	00000101	00000001
5	1: rem = rem - div	000111	00000101	11111100
	2b: rem < 0 $\Rightarrow$ +Div, shift, Q0=0	001110	00000101	00000001
	3: shift div right	001110	00000010	00000001

6 ...

but we must end at iter. 3 because we have correct values for

$Q=3$  and remainder = 11 then,

### 3.27. half precision of 16-bits



bias = 15

$$\rightarrow (-1)^n 2^{\text{exp} - \text{bias}} (1 + \text{mantissa})$$

$$-1.5625 \times 10^{-1} \rightarrow n = 1,$$

$$1.5625 \times 10^{-1} \Rightarrow \begin{array}{l} 0.15625 \times 2 = 0.3125 \quad 0 \\ 0.3125 \times 2 = 0.625 \quad 0 \\ 0.625 \times 2 = 1.25 \quad 1 \\ 0.25 \times 2 = 0.5 \quad 0 \\ 0.5 \times 2 = 1.0 \quad 1 \end{array} \Rightarrow \begin{array}{l} \text{denormalized:} \\ 0.00101 \\ \text{normalized} \\ 1.01 \times 2^{-3} \end{array}$$

$$\rightarrow (-1)^S \times 2^{\text{exp} - \text{bias}} \times (1 + \text{mant}) = (-1) \times 2^{-3} \times (1.01)$$

so,  $S = 1$ ,  $\text{exp} - \text{bias} = -3$ ,  $\text{mantissa} = 0100000000$   
 $\text{exp} - 15 = -3$ ,  
 $\text{exp} = 12 \rightarrow 1100$

$$\rightarrow \boxed{1 \mid 01100 \mid 0100000000}$$

the range is smaller than of single precision  
 bcs numbers through  $2^{-14}$  to  $2^{15}$ , whereas  
 single precision has  $2^{-126}$  to  $2^{127}$  range.

Accuracy is 10 bits while single prec. has 23 bits.

### 3.29. sum of floating point numbers

$$A = 26.125 \rightarrow 11010 + 0.125 \times 2 \quad 0.250 \quad 0 \rightarrow 11010.001_2 \quad \left. \begin{array}{l} \text{exp} = 4 + \text{bias} = 19 \\ 19 \rightarrow 10011_2 \\ \text{mant} \rightarrow 1010001 \dots \end{array} \right\}$$

$$\begin{array}{l} 0.125 \times 2 \quad 0.5 \quad 0 \\ 0.5 \times 2 \quad 1.0 \quad 1 \end{array} \rightarrow 1.1010001_2 \times 2^4$$

$$A = 0 \quad 10011 \quad 1010001000$$

$$B: 0.4150390625 \rightarrow 0110101001 \text{ using above method} \rightarrow 1.10101001 \times 2^{-2}$$

$$\rightarrow \text{exp} = -2 + \text{bias} = 13 \rightarrow 01101_2$$

$$B = 0 \quad 01101 \quad 1010100100$$

$19 - 13 = 6$  bits alignment needed  $\rightarrow$  align to higher num.  $\rightarrow$  shift B's mant. right by 6

$$B \Rightarrow 0.00000110101 \rightarrow \text{with 3 bits} \rightarrow 0.0000011010 \mid 101$$

$$\text{Add mantissas} \rightarrow \begin{array}{r} 1.1010001000 \\ + 0.0000011010 \\ \hline 1.1010010010 \quad 1000 \end{array}$$

Hence:

$$0 \quad 10011 \quad 1010100010 \quad \square$$