

## ECE 485/585 – Computer Organization and Design

### HOMEWORK #3 SOLUTION

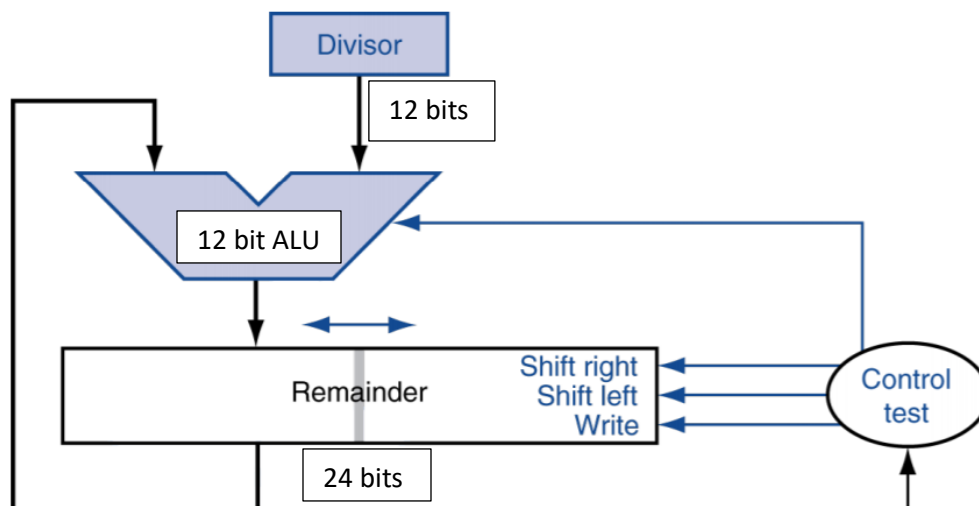
1. Provided numbers:  $6FD_{16}$  and  $273B_{16}$   
 $6FD_{16} = 0110\ 1111\ 1101\ 0100 = +28628$   
 $273B_{16} = 0010\ 0111\ 0011\ 1011 = +10043$   
 $6FD_{16} - 273B_{16} = 6FD_{16} + (-273B_{16})$   
2's comp. of  $273B = 1101\ 1000\ 1100\ 0100 = -10043$   
$$\begin{array}{r} 0110\ 1111\ 1101\ 0100 \\ + 1101\ 1000\ 1100\ 0101 \\ \hline 1\ 0100\ 1000\ 1001\ 1001 = +4899_{16} \end{array}$$
  
(pos. and neg. added together => ignore overflow)  
Answer:  $4899_{16}$
2. Provided numbers:  $174_{10}$  and  $85_{10}$   
 $174_{10} = 1010\ 1110 = -46$   
 $85_{10} = 0101\ 0101 = +85$   
2's comp of 46 =  $1101\ 0010$   
$$\begin{array}{r} 1101\ 00 \\ + 0101\ 0101 \\ \hline 1\ 0010\ 0111 = +39 \end{array}$$
  
Answer: +39, neither overflow nor underflow (ignored due to adding pos. and neg. number)
3. See the table on the next page

## ECE 485/585 – Computer Organization and Design

Iteration	Step	Quotient	Divisor	Remainder
0	Initial Values	000 000	010 110 000 000	000 000 111 100
1	1: Rem = Rem - Div	000 000	010 110 000 000	101 010 111 100
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 000	010 110 000 000	000 000 111 100
	3: Shift Div right	000 000	001 011 000 000	000 000 111 100
2	1: Rem = Rem - Div	000 000	001 011 000 000	110 101 111 100
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 000	001 011 000 000	000 000 111 100
	3: Shift Div right	000 000	000 101 100 000	000 000 111 100
3	1: Rem = Rem - Div	000 000	000 101 100 000	111 011 011 100
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 000	000 101 100 000	000 000 111 100
	3: Shift Div right	000 000	000 010 110 000	000 000 111 100
4	1: Rem = Rem - Div	000 000	000 010 110 000	111 110 001 100
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 000	000 010 110 000	000 000 111 100
	3: Shift Div right	000 000	000 001 011 000	000 000 111 100
5	1: Rem = Rem - Div	000 000	000 001 011 000	111 111 100 100
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 000	000 001 011 000	000 000 111 100
	3: Shift Div right	000 000	000 000 101 100	000 000 111 100
6	1: Rem = Rem - Div	000 000	000 000 101 100	000 000 010 000
	2a: Rem > 0 => sll Q, Q0 =1	000 001	000 000 101 100	000 000 010 000
	3: Shift Div right	000 001	000 000 010 110	000 000 010 000
7	1: Rem = Rem - Div	000 001	000 000 010 110	111 111 101 010
	2b: Rem < 0 => +Div, sll Q, Q0=0	000 010	000 000 010 110	000 000 010 000
	3: Shift Div right	000 010	000 000 001 011	000 000 010 000

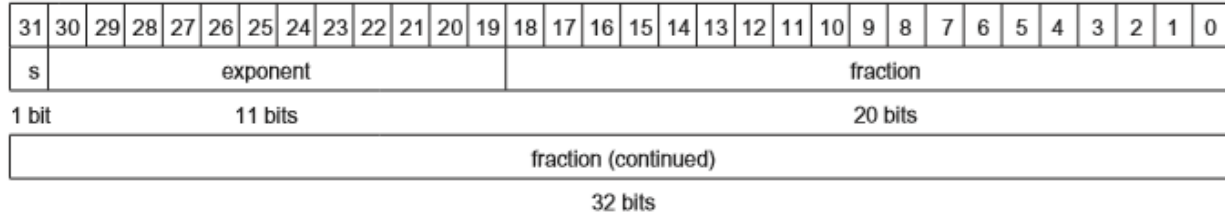
Result: Quotient = 2, Remainder = 16.  $2 \times 22 = 44 + 16 = 60$  as we expect

The figure below shows the changes to our divisor design: We need to decrease the divisor and ALU from 32 bits to 12 bits. Also, we drop the remainder from 64 bits to 24 bits.



## ECE 485/585 – Computer Organization and Design

4. Write down the binary representation of the double precision floating point number 56.93. Double precision floating points are of the following form:  
(please note that this figure from the text has an extra bit for the exponent, which should stop at bit 20)



First we need to convert 56.93 into binary

$$56 = 111000_2$$

For the fraction:

$$\begin{aligned}
 0.93 \times 2 &= 1.86 & \text{Integral Part} &= 1 \\
 0.86 \times 2 &= 1.72 & \text{Integral Part} &= 1 \\
 0.72 \times 2 &= 1.44 & \text{Integral Part} &= 1 \\
 0.44 \times 2 &= 0.88 & \text{Integral Part} &= 0 \\
 0.88 \times 2 &= 1.76 & \text{Integral Part} &= 1 \\
 0.76 \times 2 &= 1.52 & \text{Integral Part} &= 1 \\
 0.52 \times 2 &= 1.04 & \text{Integral Part} &= 1 \\
 0.04 \times 2 &= 0.08 & \text{Integral Part} &= 0 \\
 0.08 \times 2 &= 0.16 & \text{Integral Part} &= 0 \\
 0.16 \times 2 &= 0.32 & \text{Integral Part} &= 0 \\
 0.32 \times 2 &= 0.64 & \text{Integral Part} &= 0 \\
 0.64 \times 2 &= 1.28 & \text{Integral Part} &= 1 \\
 0.28 \times 2 &= 0.56 & \text{Integral Part} &= 0 \\
 0.56 \times 2 &= 1.12 & \text{Integral Part} &= 1 \\
 0.12 \times 2 &= 0.24 & \text{Integral Part} &= 0 \\
 0.24 \times 2 &= 0.48 & \text{Integral Part} &= 0 \\
 0.48 \times 2 &= 0.96 & \text{Integral Part} &= 0 \\
 0.96 \times 2 &= 1.92 & \text{Integral Part} &= 1 \\
 0.92 \times 2 &= 1.84 & \text{Integral Part} &= 1 \\
 0.84 \times 2 &= 1.68 & \text{Integral Part} &= 1 \\
 0.68 \times 2 &= 1.36 & \text{Integral Part} &= 1 \\
 0.36 \times 2 &= 0.72 & \text{Integral Part} &= 0 \\
 0.72 \times 2 &= 1.44 & \text{Integral Part} &= 1
 \end{aligned}$$

\*\*Found same number as iteration 3, so will repeat endlessly

So in total:  $56.93 = 111000.1110111000010100011110101110000101000111101\dots$

First we shift our number so that we have a leading 1 and then the decimal place:

$$1.11000111011100001010001111010\dots \times 2^5$$

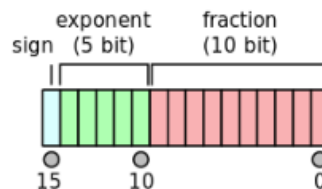
We have a sign bit of 0 and now need to calculate the exponent =  $5 + 1023 = 1028$

$$1028 = 100\ 0000\ 0100_2$$

We have a sign bit of 0 in this case.

3	3	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1													
1	0	9	8	7	6	5	4	3	2	1	0	9	8	7	6	5	4	3	2	1	0	9	8	7	6	5	4	3	2	1	0		
S	exponent											Fraction																					
0	100 0000 0100											1100 0111 0111 0000 1010																					
Fraction (continued)																																	
0011 1101 0111 0000 1010 0011 1101 0111																																	

5. Write down the half-precision representation of the following number:  $-1.585 \times 10^{-1}$ . Half Precision floating point numbers are of the following form:



0.1585\*2 = 0.317 Integral Part = 0  
0.317 \*2 = 0.634 Integral Part = 0  
0.634 \*2 = 1.268 Integral Part = 1  
0.268 \*2 = 0.536 Integral Part = 0  
0.536 \*2 = 1.072 Integral Part = 1  
0.072 \*2 = 0.144 Integral Part = 0  
0.144 \*2 = 0.288 Integral Part = 0  
0.288 \*2 = 0.576 Integral Part = 0  
0.576 \*2 = 1.152 Integral Part = 1  
0.152 \*2 = 0.304 Integral Part = 0  
0.304 \*2 = 0.608 Integral Part = 0  
0.608 \*2 = 1.216 Integral Part = 1  
0.216 \*2 = 0.432 Integral Part = 0

This results in  $0.0010100010010 = 1.0100010010 \cdot 2^{-3}$

$$12 = 01100_2$$

So we get the following floating point number:

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
S	exponent					Fraction									
1	011 00					01 0001 0010									

When compared to single-precision floating point, we receive a better range for the number of bits that we use (5/16 compared to 8/32) however we sacrifice precision in doing so.