#### Lecture 10: Division, Floating Point

- Today's topics:
  - Division

IEEE 754 representations - Floating Point

# Divide Example

• Divide  $7_{ten}$  (0000 0111<sub>two</sub>) by  $2_{ten}$  (0010<sub>two</sub>)

lter	Step	Quot	Divisor	Remainder
0	Initial values	0000	0010 0000	0000 0111
1	Rem = Rem – Div	0000	0010 0000	1110 0111
	Rem < 0 🗲 +Div, shift 0 into Q	0000	0010 0000	0000 0111
	Shift Div right	0000	0001 0000	0000 0111
2	Same steps as 1	0000	0001 0000	1111 0111
		0000	0001 0000	0000 0111
		0000	0000 1000	0000 0111
3	Same steps as 1	0000	0000 0100	0000 0111
4	Rem = Rem – Div	0000	0000 0100	0000 0011
	Rem >= 0 ➔ shift 1 into Q	0001	0000 0100	0000 0011
//	Shift Div right	0001	0000 0010	0000 0011
5	Same steps as 4	0011	0000 0001	0000 0001

# Hardware for Division



A comparison requires a subtract; the sign of the result is examined; if the result is negative, the divisor must be added back

Similar to multiply, results are placed in Hi (remainder) and Lo (quotient)

## **Efficient Division**



## **Divisions involving Negatives**

- Simplest solution: convert to positive and adjust sign later
- Note that multiple solutions exist for the equation:
  Dividend = Quotient x Divisor + Remainder

+7	div +2	Quo =	Rem =
-7	div +2	Quo =	Rem =
+7	div -2	Quo =	Rem =
-7	div -2	Quo =	Rem =

## **Divisions involving Negatives**

- Simplest solution: convert to positive and adjust sign later
- Note that multiple solutions exist for the equation:
  Dividend = Quotient x Divisor + Remainder

+7 div +2	Quo = +3	Rem = +1	Que = 4	Ren = -1
-7 div +2	Quo = -3	Rem = -1		
+7 div -2	Quo = -3	Rem = +1		
-7 div -2	Quo = +3	Rem = -1	anor 4	Rem = +1

Convention: Dividend and remainder have the same sign Quotient is negative if signs disagree These rules fulfil the equation above

#### Take Homes

Grade school algorithms are commonly used – the algorithms are even easier in binary (mult by 1 and 0)

They can be implemented in hardware with shifts, add, sub, checks

• To improve efficiency, look for ineffectuals – are only some bits changing in every step – allows us to use narrow adders and registers – allows us to pack more operands in one register

• Can also improve speed by throwing more transistors and parallel computations at the problem Wallace T(ee)







- This is the biased notation, where a bias is subtracted from the exponent field to yield the true exponent
- IEEE 754 single-precision uses a bias of 127 (since the exponent must have values between -127 and 128)...double precision uses a bias of 1023

Final representation:  $(-1)^{S} \times (1 + Fraction) \times 2^{(Exponent - Bias)}$ 

## Sign and Magnitude Representation



- Largest number that can be represented: 2.0 x 2<sup>128</sup> = 2.0 x 10<sup>38</sup> (not really – see upcoming details)
- Smallest number that can be represented: 1.0 x 2<sup>-127</sup> = 2.0 x 10<sup>-38</sup> (not really – see upcoming details)
- Overflow: when representing a number larger than the max; Underflow: when representing a number smaller than the min
- Double precision format: occupies two 32-bit registers:



11

#### Details

- The number "0" has a special code so that the implicit 1 does not get added: the code is all 0s (it may seem that this takes up the representation for 1.0, but given how the exponent is represented, that's not the case) (see discussion of denorms in the textbook)
- The largest exponent value (with zero fraction) represents +/- infinity
- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of 0/0 or (infinity minus infinity)
- Note that these choices impact the smallest and largest numbers that can be represented





#### **Examples**

Final representation: (-1)<sup>S</sup> x (1 + Fraction) x 2<sup>(Exponent – Bias)</sup>

- Represent -0.75<sub>ten</sub> in single and double-precision formats
  Single: (1 + 8 + 23)
  1 0111 1110 1000...000
  -++27
  Double: (1 + 11 + 52)
  1 0111 1111 110 1000...000
- What decimal number is represented by the following single-precision number?
  - $1 \quad 1000 \ 0001 \quad 01000...0000$

-5.0

Final representation: (-1)<sup>S</sup> x (1 + Fraction) x 2<sup>(Exponent – Bias)</sup>

• Represent 36.90625<sub>ten</sub> in single-precision format



Final representation: (-1)<sup>S</sup> x (1 + Fraction) x 2<sup>(Exponent – Bias)</sup>

We've calculated that  $36.90625_{ten} = 100100.1110100...0$  in binary Normalized form = 1.001001110100...0 x  $2^5$ (had to shift 5 places to get only one bit left of the point)

The sign bit is 0 (positive number) The fraction field is 001001110100...0 (the 23 bits after the point) The exponent field is 5 + 127 (have to add the bias) = 132, which in binary is 10000100

The IEEE 754 format is 0 10000100 001001110100.....0 sign exponent 23 fraction bits