

FLOATING POINT

Mahdi Nazm Bojnordi

Assistant Professor

School of Computing

University of Utah

Overview

- This lecture
 - ▣ Hardware for division
 - ▣ MIPS instructions for division
 - ▣ Signed division
 - ▣ Floating point numbers

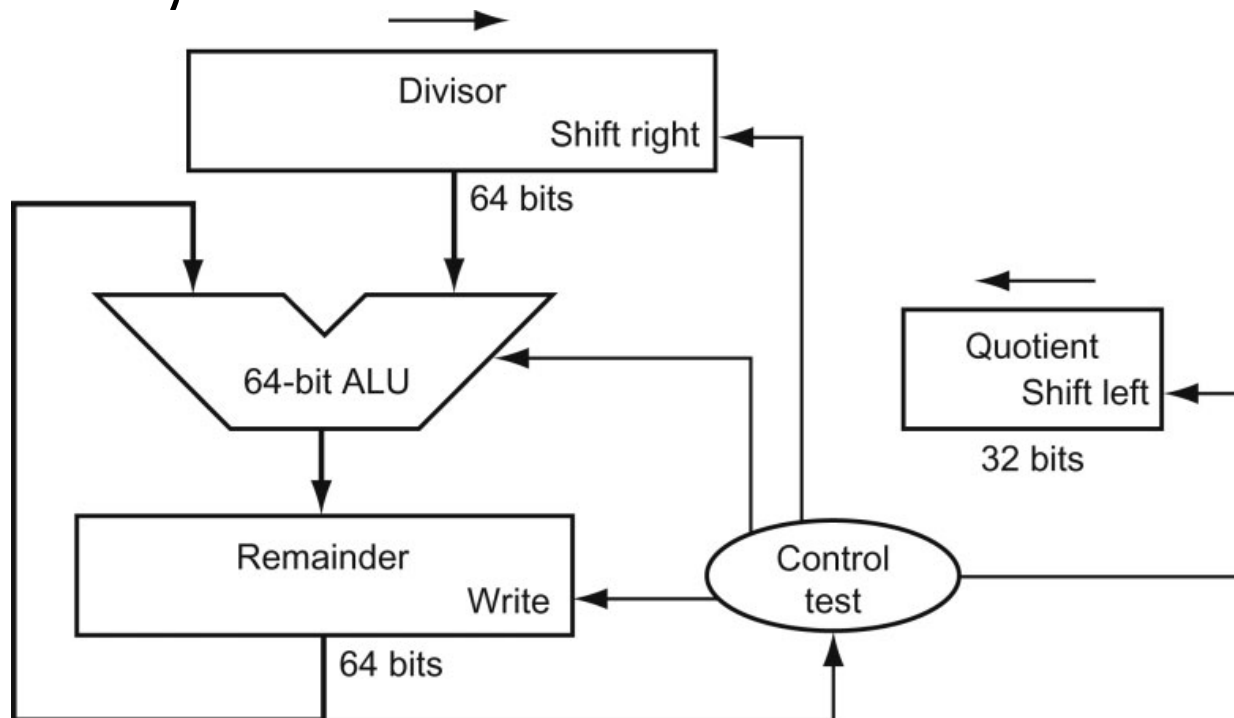
Recall: Division Example

□ Divide 7_{ten} (0000 0111_{two}) by 2_{ten} (0010_{two})

Iter	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 (pp. 189-193)

- 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)



MIPS Instructions for Division

- MIPS ignores overflow in division
- Signed division (div)

`div` `$s2, $s3` computes the division and stores
it in two “internal” registers that
can be referred to as `hi` and `lo`

`mfhi` `$s0` moves the remainder into `$s0`
`mflo` `$s1` moves the quotient into `$s1`

- Unsigned division (divu)

`divu` `$s2, $s3`

`mfhi` `$s0`
`mflo` `$s1`

Signed Division

- **Simplest solution:** convert to positive and adjust sign later
- Note that multiple solutions exist for the equation
 - ▣ $\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}$

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

Signed Division

- **Simplest solution:** convert to positive and adjust sign later
- Note that multiple solutions exist for the equation
 - ▣ $\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}$

+7	div	+2	Quo = +3	Rem = +1
-7	div	+2	Quo = -3	Rem = -1
+7	div	-2	Quo = -3	Rem = +1
-7	div	-2	Quo = +3	Rem = -1

Signed Division

- **Simplest solution:** convert to positive and adjust sign later
- Note that multiple solutions exist for the equation
 - ▣ $\text{Dividend} = \text{Quotient} \times \text{Divisor} + \text{Remainder}$

+7	div	+2	Quo = +3	Rem = +1
-7	div	+2	Quo = -3	Rem = -1
+7	div	-2	Quo = -3	Rem = +1
-7	div	-2	Quo = +3	Rem = -1

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

Floating Point Numbers

□ Example: how to represent 32.5

▣ weights smaller than 1

100000.100_{two}

Floating Point Numbers

- Example: how to represent 32.5

- ▣ weights smaller than 1

100000.100_{two}

- Normalized scientific notation: leave a single non-zero digit to the left of the point

- ▣ 3.25×10^1

$$1.00000100_{\text{two}} * 2^5 = (1 * 2^0 + 1 * 2^{-6}) * 2^5$$

- The IEEE 754 standard

Sign-Magnitude Representation

- Since we are only representing normalized numbers, we are guaranteed that the number is of the form $1.xxxx$.
- Every 32-bit number has three fields: sign (S), exponent (E), and fraction (F)

▣ $\text{value} = (-1)^S \times (1 + F) \times 2^{E-127}$



Sign-Magnitude Representation

- Since we are only representing normalized numbers, we are guaranteed that the number is of the form $1.xxxx$.

- Every 32-bit number has three fields: sign (S), exponent (E), and fraction (F)

$$\blacksquare \text{ value} = (-1)^S \times (1 + F) \times 2^{E-127}$$



- Example: 32.5

$$\blacksquare 32.5 = 100000.100 = 1.00000100_{\text{two}} \times 2^5$$

Sign-Magnitude Representation

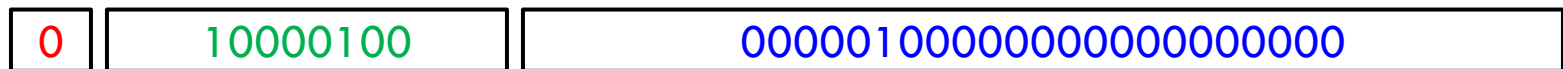
- Since we are only representing normalized numbers, we are guaranteed that the number is of the form $1.xxxx$.
- Every 32-bit number has three fields: sign (S), exponent (E), and fraction (F)

$$\blacksquare \text{ value} = (-1)^S \times (1 + F) \times 2^{E-127}$$



- Example: 32.5

$$\blacksquare 32.5 = 100000.100 = 1.00000100_{\text{two}} \times 2^5$$



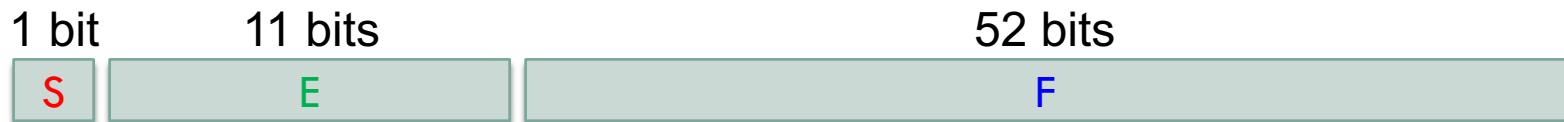
Sign-Magnitude Representation

- Since we are only representing normalized numbers, we are guaranteed that the number is of the form $1.xxxx$.
- Every 32-bit number has three fields: sign (S), exponent (E), and fraction (F)

$$\blacksquare \text{ value} = (-1)^S \times (1 + F) \times 2^{E-127}$$



- Using more bits
 - ▣ Increase E bits to represent a wider range of numbers
 - ▣ Increase F bits to represent more precision
- Double precision format with 64 bits



Single Precision Floating Point



- ☐ Largest number that can be represented
- ☐ Smallest number that can be represented

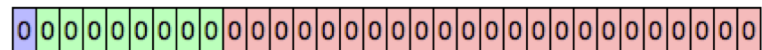
Single Precision Floating Point



- Largest number that can be represented
 - ▣ $2.0 \times 2^{128} = 2.0 \times 10^{38}$
- Smallest number that can be represented
 - ▣ $1.0 \times 2^{-127} = 2.0 \times 10^{-38}$
- Overflow
 - ▣ representing a number larger than the one above
- Underflow
 - ▣ representing a number smaller than the one above

Special Notes

- The number “0” has a special code so that the implicit 1 does not get added
 - ▣ see discussion of denorms (pg. 222) in the textbook

 = 0

$$1 \times 2^{-126} \times 0$$

- 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)

Special Notes

- The number “0” has a special code so that the implicit 1 does not get added
 - ▣ see discussion of denorms (pg. 222) in the textbook

[illegible]

$$1 \times 2^{-126} \times 0$$

- The largest exponent value (with zero fraction) represents \pm infinity

1 1 1 1 1 1 1 1 0 = -Infinity

$$-1 \times 2^{128} \times 1$$

- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of $0/0$ or (infinity minus infinity)

Special Notes

- The number “0” has a special code so that the implicit 1 does not get added
 - ▣ see discussion of denorms (pg. 222) in the textbook

[illegible]

$$1 \times 2^{-126} \times 0$$

- The largest exponent value (with zero fraction) represents +/- infinity

[illegible]

$$-1 \times 2^{128} \times 1$$

- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of $0/0$ or (infinity minus infinity)

NaN

$$-1 \times 2^{128} \times 1.00000001$$

Special Notes

- The number “0” has a special code so that the implicit 1 does not get added
 - ▣ see discussion of denorms (pg. 222) in the textbook

[illegible]

$1 \times 2^{-126} \times 0$

- The largest exponent value (with zero fraction) represents $+\infty$ or $-\infty$

[illegible]

$$-1 \times 2^{128} \times 1$$

- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of $0/0$ or (infinity minus infinity)

[illegible]

$$-1 \times 2^{128} \times 1.00000001$$

These choices impact the smallest and largest numbers that can be represented

Example 1

- Represent -0.75_{ten} in single- and double-precision formats

$$\text{value} = (-1)^{\textcolor{red}{S}} \times (1 + \textcolor{blue}{F}) \times 2^{\textcolor{green}{E} - 127}$$

- Single precision

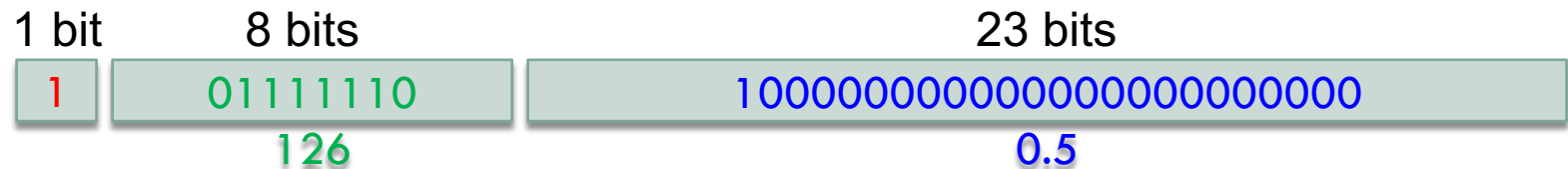


Example 1

- Represent -0.75_{ten} in single- and double-precision formats

$$-0.75 = -1 \times 0.11_{\text{two}} = -1 \times 1.1_{\text{two}} \times 2^{-1} = (-1)^{\textcolor{red}{1}} \times (1 + \textcolor{blue}{F}) \times 2^{\textcolor{green}{(E - 127)}}$$

- Single precision



- Represent -0.75_{ten} in single- and double-precision formats

- Single precision



[illegible]

Example 2

- Represent 3.40625_{ten} in single- and double-precision formats

$$\text{value} = (-1)^{\textcolor{red}{S}} \times (1 + \textcolor{blue}{F}) \times 2^{\textcolor{green}{E} - 127}$$

- Single precision

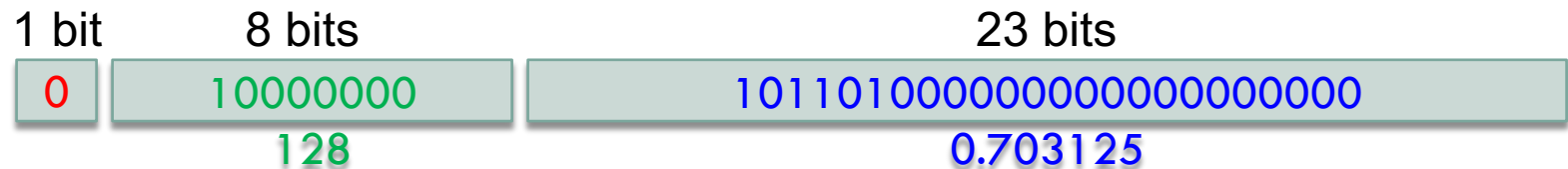


Example 2

- Represent 3.40625_{ten} in single- and double-precision formats

$$3.40625 = 11.01101_{\text{two}} = 1.101101_{\text{two}} \times 2^1 = (-1)^0 \times (1 + F) \times 2^{(E - 127)}$$

- Single precision

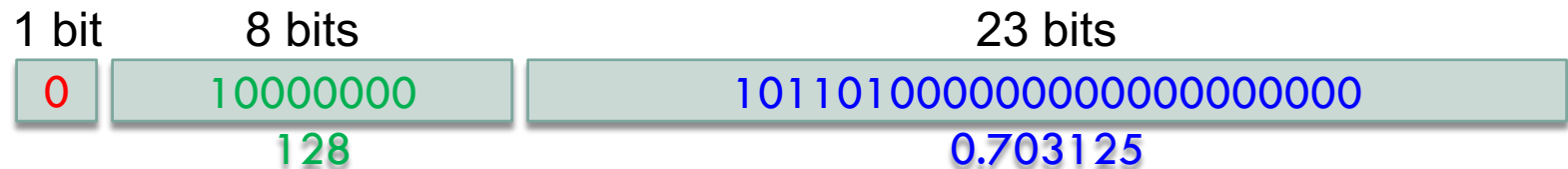


Example 2

- Represent 3.40625_{ten} in single- and double-precision formats

$$3.40625 = 11.01101_{\text{two}} = 1.101101_{\text{two}} \times 2^1 = (-1)^0 \times (1 + F) \times 2^{(E - 127)}$$

- Single precision



- Double precision

$$1.101101_{\text{two}} \times 2^1 = (-1)^0 \times (1 + F) \times 2^{(E - 1023)}$$

