

## Chapter Three

## Numbers

- Bits are just bits (no inherent meaning)
  - conventions define relationship between bits and numbers
- Binary numbers (base 2)
  - 0000 0001 0010 0011 0100 0101 0110 0111 1000 1001...
  - decimal:  $0 \dots 2^n - 1$
- Of course it gets more complicated:
  - numbers are finite (overflow)
  - fractions and real numbers
  - negative numbers
  - e.g., no MIPS subi instruction; addi can add a negative number
- How do we represent negative numbers?
  - i.e., which bit patterns will represent which numbers?

Handwritten notes:  $18_{10}$ ,  $11_2$ ,  $1_9$ ,  $9$ ,  $11$  with vertical lines and arrows indicating binary conversion steps.

Handwritten note: 885

## Possible Representations

Sign Magnitude	One's Complement	Two's Complement
000 = +0	000 = +0	000 = +0
001 = +1	001 = +1	001 = +1
010 = +2	010 = +2	010 = +2
011 = +3	011 = +3	011 = +3
100 = -0	100 = -3	100 = -4
101 = -1	101 = -2	101 = -3
110 = -2	110 = -1	110 = -2
111 = -3	111 = -0	111 = -1

Handwritten note:  $0 - (-4)$

- Issues: balance, number of zeros, ease of operations
- Which one is best? Why?

Handwritten equation:  $X + (-x) = 000$

## MIPS

- 32 bit signed numbers:

0000 0000 0000 0000 0000 0000 0000 0000	= $0_{\text{ten}}$
0000 0000 0000 0000 0000 0000 0000 0001	= $+1_{\text{ten}}$
0000 0000 0000 0000 0000 0000 0000 0010	= $+2_{\text{ten}}$
...	
0111 1111 1111 1111 1111 1111 1111 1110	= $+2,147,483,646_{\text{ten}}$
0111 1111 1111 1111 1111 1111 1111 1111	= $+2,147,483,647_{\text{ten}}$
1000 0000 0000 0000 0000 0000 0000 0000	= $-2,147,483,648_{\text{ten}}$
1000 0000 0000 0000 0000 0000 0000 0001	= $-2,147,483,647_{\text{ten}}$
1000 0000 0000 0000 0000 0000 0000 0010	= $-2,147,483,646_{\text{ten}}$
...	
1111 1111 1111 1111 1111 1111 1111 1101	= $-3_{\text{ten}}$
1111 1111 1111 1111 1111 1111 1111 1110	= $-2_{\text{ten}}$
1111 1111 1111 1111 1111 1111 1111 1111	= $-1_{\text{ten}}$

Handwritten labels: maxint, minint

## Two's Complement Operations

- Negating a two's complement number: invert all bits and add 1
  - remember: "negate" and "invert" are quite different!
- Converting n bit numbers into numbers with more than n bits:
  - MIPS 16 bit immediate gets converted to 32 bits for arithmetic
  - copy the most significant bit (the sign bit) into the other bits
  - 0010 -> 0000 0010
  - 1010 -> 1111 1010
- "sign extension" (lbu vs. lb)

## Addition & Subtraction

Just like in grade school (carry/borrow 1s)

$$\begin{array}{r} 0111 \quad 0111 \quad 0110 \\ + 0110 \quad - 0110 \quad - 0101 \end{array}$$

Two's complement operations easy

- subtraction using addition of negative numbers

$$\begin{array}{r} 0111 \\ + 1010 \end{array}$$

- Overflow (result too large for finite computer word):
    - e.g., adding two n-bit numbers does not yield an n-bit number
- $$\begin{array}{r} 0111 \\ + 0001 \\ - 1000 \end{array}$$
- note that overflow term is somewhat misleading, it does not mean a carry "overflowed"

Overflow

+	+	+	S	N	+	-	+
+	+	-	N	S	+	-	-
-	+	+	N	S	-	-	+
-	+	-	S	N	-	-	-

$x - (y)$

### Detecting Overflow

- No overflow when adding a positive and a negative number
- No overflow when signs are the same for subtraction
- Overflow occurs when the value affects the sign:
  - overflow when adding two positives yields a negative
  - or, adding two negatives gives a positive
  - or, subtract a negative from a positive and get a negative
  - or, subtract a positive from a negative and get a positive
- Consider the operations  $A + B$ , and  $A - B$ 
  - Can overflow occur if  $B$  is 0?
  - Can overflow occur if  $A$  is 0?

©2004 Morgan Kaufmann Publishers 7

### Effects of Overflow

- An exception (interrupt) occurs
  - Control jumps to predefined address for exception
  - Interrupted address is saved for possible resumption
- Details based on software system / language
  - example: flight control vs. homework assignment
- Don't always want to detect overflow
  - new MIPS instructions: **addu, addiu, subu**

*note: addiu still sign-extends!*  
*note: sltu, sltiu for unsigned comparisons*

©2004 Morgan Kaufmann Publishers 8

### Multiplication

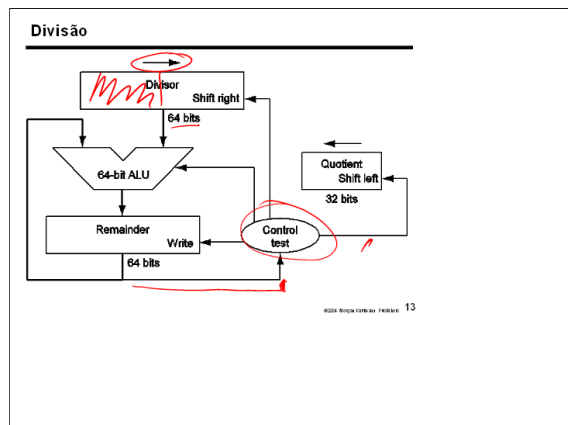
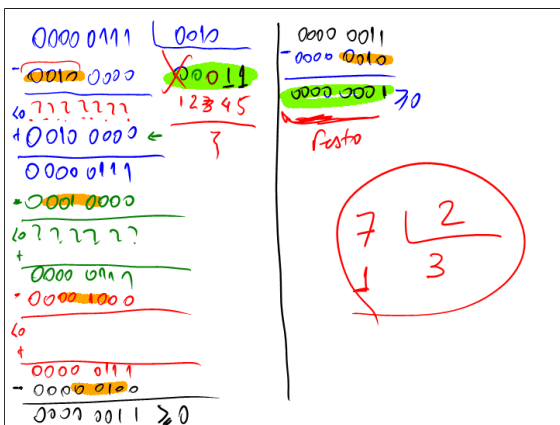
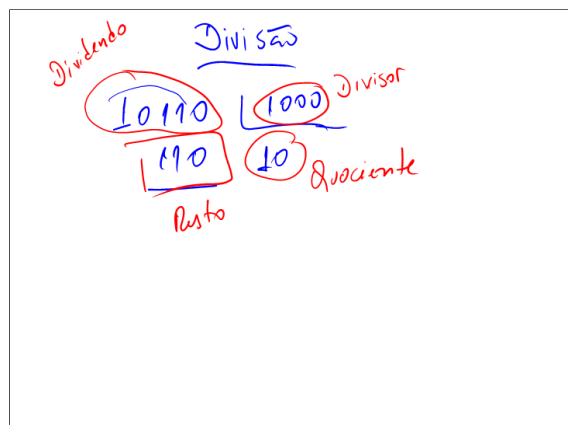
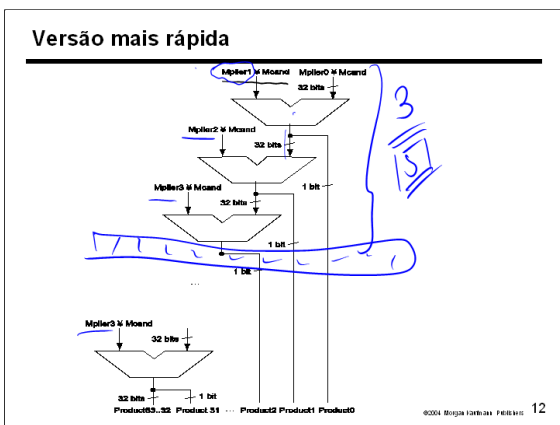
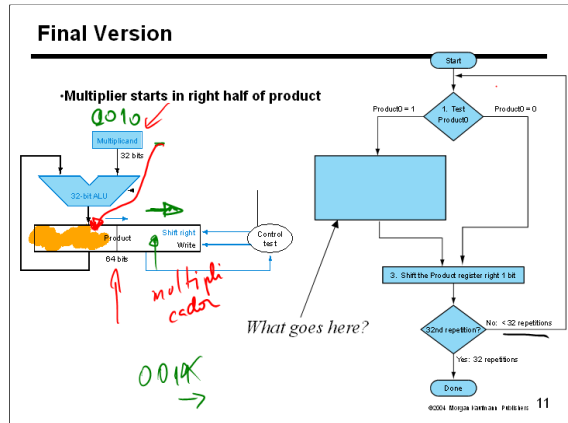
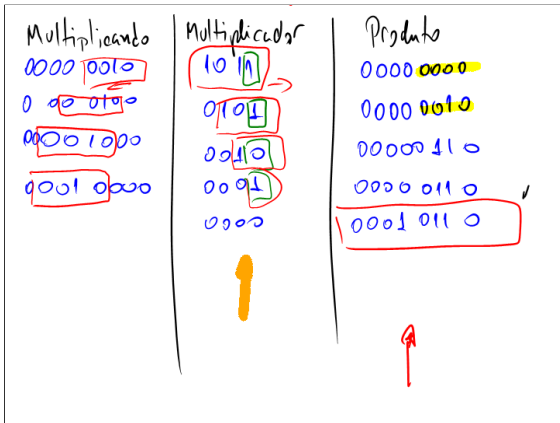
- More complicated than addition
  - accomplished via shifting and addition
- More time and more area
- Let's look at 3 versions based on a gradeschool algorithm
 

$$\begin{array}{r} 0010 \text{ (multiplicand)} \\ \times 1011 \text{ (multiplier)} \\ \hline \end{array}$$
- Negative numbers: convert and multiply
  - there are better techniques, we won't look at them

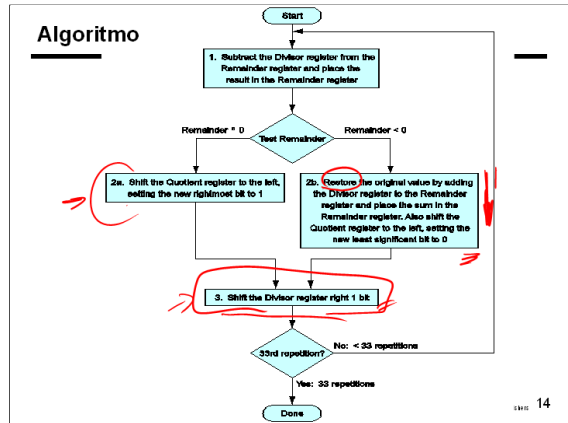
©2004 Morgan Kaufmann Publishers 9

### Multiplication: Implementation

©2004 Morgan Kaufmann Publishers 10



Quociente	Divisor	Resto
0000	0010 0000	0000 0111
0000	0001 0000	0000 0111
0000	0000 1000	0000 0111
0000	0000 0100	0000 0111
0001	0000 0010	0000 0001
0011		0000 0000



int i, j      4.75  
 float a, b      π

134.37  
 $1343.7 \times 10^{-1}$   
 $1.3437 \times 10^2$  \*

3,333333...

### Floating Point

- We need a way to represent
  - numbers with fractions, e.g., 3.1416
  - very small numbers, e.g., 0.000000001
  - very large numbers, e.g.,  $3.15576 \times 10^9$
- Representation:
  - sign, exponent, significand:  $(-1)^{\text{sign}} \times \text{significand} \times 2^{\text{exponent}}$
  - more bits for significand gives more accuracy
  - more bits for exponent increases range
- Overflow
- Underflow

©2004 Morgan Kaufmann Publishers 15

~~1.11~~ × 2<sup>2</sup>

### Como representar? IEEE 754

- Números normalizados
  - Números da forma 1.xxxxxxxx
  - Não é necessário armazenar o 1.
- Representação com 32 bits (precisão simples)

3	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0
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	Expoente 8										Mantissa 23																					

- Representação com 64 bits (precisão dupla)

3	2	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0			
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	Expoente 11										Mantissa 52																									

3	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
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		
Mantissa																																	

©2004 Morgan Kaufmann Publishers 16

### IEEE 754 floating-point standard

- IEEE 754 floating point standard:
  - single precision: 8 bit exponent, 23 bit significand
  - double precision: 11 bit exponent, 52 bit significand
- Exponent is "biased" to make sorting easier
  - all 0s is smallest exponent all 1s is largest
  - bias of 127 for single precision and 1023 for double precision
  - summary:  $(-1)^{\text{sign}} \times (1 + \text{significand}) \times 2^{\text{exponent} - \text{bias}}$
- Example:
  - decimal:  $-0.75 = -(1/2 + 1/4)$
  - binary:  $-0.11 = -1.1 \times 2^{-1}$
  - floating point: exponent = 126 = 01111110
  - IEEE single precision: 10111111010000000000000000000000

©2004 Morgan Kaufmann Publishers 17

Handwritten calculation showing the addition of two floating-point numbers:

$$1.610 \times 10^{-1} + 1.781 \times 10^2$$

The second number is shifted to match the exponent of the first:

$$1.781 \times 10^2 = 1781.0 \times 10^{-1}$$

$$1.610 \times 10^{-1} + 1781.0 \times 10^{-1} = 1782.610 \times 10^{-1}$$

The final result is  $1.782610 \times 10^2$ .

Handwritten calculation showing the addition of two floating-point numbers in binary:

$$0.11 = 1.1 \times 2^{-1} + 1.1 \times 2^{-4}$$

$$0.00011 \times 2^4 = 110000.0 \times 2^{-1}$$

The exponents are aligned:

$$126 + 131 = 257$$

$$126 + 131 - 127 = 130$$

The result is  $1.100011 \times 2^{-1}$ .

### Exemplos

- Comparar em binário
  - 0.75; -0.75; 4.25; 0.625; 19

©2004 Morgan Kaufmann Publishers 18

### Constantes

Precisão Simples		Precisão Dupla		Valor
Expoente	Mantissa	Expoente	Mantissa	
0	0	0	0	0
0	<> 0	0	<> 0	Número não normalizado
1-254	Qualquer	1-2046	Qualquer	Número em ponto flutuante
255	0	2047	0	Infinito
255	<> 0	2047	<> 0	NaN (Not a Number)

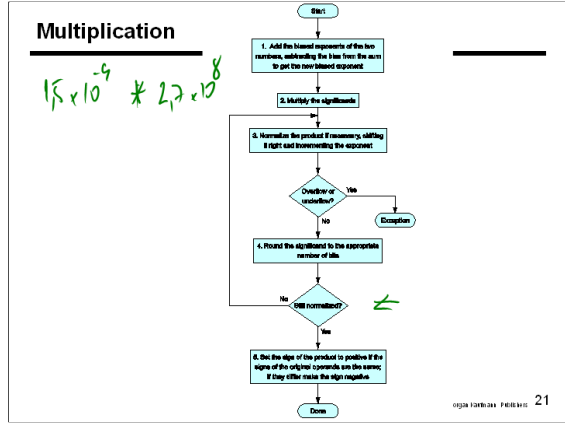
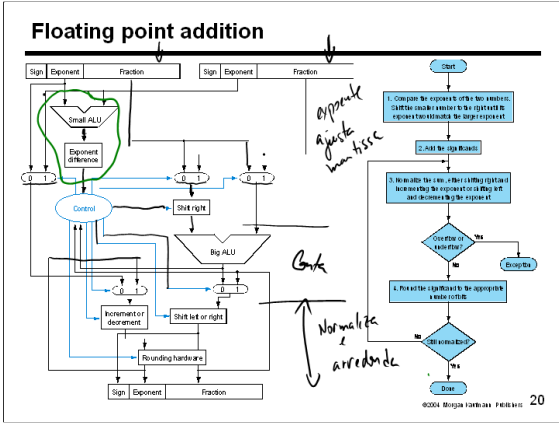
©2004 Morgan Kaufmann Publishers 19

Handwritten diagram of a floating-point number format:

$f$  |  $exp$  | mantissa

$1 \times 2^{-1}$

$1$



### Precisão

- $X + 1 - X = 1?$
- Internamente, o processador armazena os números de ponto flutuante com ao menos 2 bits a mais: *guard e round*
- Um terceiro bit, o *sticky* indica se algum conteúdo significativo foi perdido em arredondamento
- Formas de arredondar:
  - Em direção a +infinito
  - Em direção a -infinito
  - Truncar
  - Em direção ao par mais próximo

Handwritten notes: 'float', 'int', '2<sup>32</sup>', '2<sup>32</sup>'

©2004 Morgan Kaufmann Publishers 22

Handwritten notes:  $(X) + 1 - X$ ,  $X \gg \gg \gg 1$

### Floating Point Complexities

- Operations are somewhat more complicated (see text)
- In addition to overflow we can have "underflow"
- Accuracy can be a big problem
  - IEEE 754 keeps two extra bits, guard and round
  - four rounding modes
  - positive divided by zero yields "infinity"
  - zero divide by zero yields "not a number"
  - other complexities
- Implementing the standard can be tricky
- Not using the standard can be even worse
  - see text for description of 80x86 and Pentium bug!

©2004 Morgan Kaufmann Publishers 23

### Chapter Three Summary

- Computer arithmetic is constrained by limited precision
- Bit patterns have no inherent meaning but standards do exist
  - two's complement
  - IEEE 754 floating point
- Computer instructions determine "meaning" of the bit patterns
- Performance and accuracy are important so there are many complexities in real machines
- Algorithm choice is important and may lead to hardware optimizations for both space and time (e.g., multiplication)
- You may want to look back (Section 3.10 is great reading!)

©2004 Morgan Kaufmann Publishers 24