Overview

COMP 3221

Microprocessors and Embedded Systems

Lectures 19 : Floating Point Number Representation – I http://www.cse.unsw.edu.au/~cs3221

September, 2003

Saeid Nooshabadi

Saeid@unsw.edu.au

COMP3221 lec19-fp-l.1

Review of Numbers

°Computers are made to deal with numbers

°What can we represent in N bits?

- Unsigned integers:
 - 0 to 2^N 1
- Signed Integers (Two's Complement)

_2(N-1) to

to 2^(N-1) - 1

°Floating Point Numbers

- ^o Motivation: Decimal Scientific Notation
 - Binary Scientific Notation
- °Floating Point Representation inside computer (binary)
 - Greater range, precision

°IEEE-754 Standard

COMP3221 lec19-fp-l.2

Saeid Nooshabadi

Other Numbers

 ^o What about other numbers?
 • Very large numbers? (seconds/century) 3,155,760,000₁₀ (3.15576₁₀ x 10⁹)
 • Very small numbers? (atomic diameter) 0.00000001₁₀ (1.0₁₀ x 10⁻⁸)
 • Rationals (repeating pattern) 2/3 (0.6666666666...)
 • Irrationals

- rationals 2^{1/2} (1.414213562373...)
- Transcendentals e (2.718...), π (3.141...)

°All represented in scientific notation

COMP3221 lec19-fp-l.4

Saeid Nooshabadi



Floating Point Number Range vs Precision (#3/4)

°In which representation can the most different values be represented?

°All can represent the same number of different values.

°Two of the systems can represent an equal number of values, more than the third.

 $N = d_1 \times 10^{d_2 d_3}$

 $N = d_1 d_2 \times 10^{d_3}$

 $N = d_1 d_2 \times 100^{d_3}$

°#3 can represent more values than the other two.

°#2 can represent more values than the other two.

°#1 can represent more values than the other two Saeid Nooshabadi

COMP3221 lec19-fp-l.9

Scientific Notation for Binary Numbers

exponent significand 1.0_{two} x 2⁻ radix (base) "binary point"

^oComputer arithmetic that supports it called <u>floating point</u>, because it represents numbers where binary point is not fixed, as it is for integers

• Declare such variable in C as float

Floating Point Number Range vs Precision (#4/4)



COMP3221 lec19-fp-l.10

Properties of a good FP Number rep.

^oRepresents many useful numbers • most of the 2^N possible are useful • How many LARGE? How many small? [°]Easy to do arithmetic (+, -, *, /) $^{\circ}$ Easy to do comparison (==, <, >) ^oNice mathematical properties

• A != B => A - B != 0

Floating Point Representation (#1/2)

[°]Normal format: +1.xxxxxxxxx_{two}*2^{yyyy}two
 [°]Multiple of Word Size (32 bits)

. 31.30 23	22 0								
S Exponent	Significand								
1 bit 8 bits °S represent	23 bits s <mark>Sign</mark>								
[°] Exponent represents y's									
Significand represents x's									
° Leading 1 in Significand is implied									
^o Represent numbers as small as 2.0 x 10 ⁻³⁸ to as large as 2.0 x 10 ³⁸									

Floating Point Representation (#2/2)

°What if result too large? (> 2.0x10 ³⁸) • <u>Overflow</u> !
 Overflow => Exponent larger than represented in 8-bit Exponent field
 What if result too small? (<0, < 2.0x10⁻³⁸) <u>Underflow!</u>
 Underflow => Negative exponent larger than represented in 8-bit Exponent field
°How to reduce chances of overflow or underflow?
COMP3221 lec19-fp-l.14 Saeid Nooshabadi

Double Precision FI. Pt. Representation

°Next Multiple of Word Size (64 bits)

3 <u>1 30</u>	20) 19	0	
S	Exponent	Significand		
1 bit	11 bits	20 bits		
	5	Significand (cont'd)		

32 bits °<u>Double Precision</u> (vs. <u>Single Precision</u>)

- C variable declared as double
- Represent numbers almost as small as 2.0 x 10⁻³⁰⁸ to almost as large as 2.0 x 10³⁰⁸
- But primary advantage is greater accuracy due to larger significand

computation using a special coprocessor.

works under the processor supervision

Fl. Pt. Hardware

^o Microprocessors do floating point

- Has its own set of registers
- [°]Most low end processors do not have Ft. Pt. Coprocessors
 - Ft. Pt. Computation by software emulation
 - ARM processor on DSLMU board does not have Ft. Pt. Coprocessor
- Some high end ARM processors do

Saeid Nooshabadi

IEEE 754 Floating Point Standard (#2/6) IEEE 754 Floating Point Standard (#1/6) ^o Single Precision, (DP similar) ^o Kahan^{*} wanted FP numbers to be used even if no FP hardware; e.g., sort records with FP ° Sign bit: 1 means negative numbers using integer compares 0 means positive ° Wanted Compare to be faster, by means of a ^o Significand: single compare operation, used for integer • To pack more bits, leading 1 implicit for numbers, especially if positive FP numbers normalized numbers. (Hidden Bit) ^o How to order 3 parts (Sign, Significand and Exponent) to simplify compare? •1 + 23 bits single, 1 + 52 bits double always true: 0 < Significand < 1 The Author of IEEE 754 FP Standard (for normalized numbers) http://www.cs.berkeley.edu/~wkahan/ ieee754status/754story.html ^oWhat about Zero? **Next Lecture** Saeid Nooshabadi Saeid Nooshabadi COMP3221 lec19-fp-l.17 COMP3221 lec19-fp-l.18 IEEE 754 Floating Point Standard (#3/6) **IEEE 754 Floating Point Standard (#4/6)** ^oWant compare FI.Pt. numbers as if integers, to help in sort ^o How to order 3 Fields in a Word? +1.xxxxxxxxx_{two}*2^{yyyy}two ° "Natural": Sign, Fraction, Exponent? Sign first part of number • Problem: If want to sort using integer • Exponent next, so big exponent => bigger No. compare operations, won't work: • 1.0 x 2²⁰ vs. 1.1 x 2¹⁰; latter looks bigger! ^oNegative Exponent? 0 10000 10100 0 11000 01010 • 2's comp? 1.0 x 2⁻¹ vs 1.0 x2⁺¹ (1/2 vs 2) ^o Exponent, Sign, Fraction? 1/2 0 1111 1111 000 0000 0000 0000 0000 0000 • Need to get sign first, 0000 0001 000 0000 0000 0000 0000 0000 0 since negative < positive 2 ^o Therefore order is Sign Exponent Fraction 1.0 x 10²⁰ > 1.1 x 10¹⁰ This notation using integer compare of 1/2 vs 2 makes 1/2 > 2!Exponent Significand COMP3221 lec COMP3221 lec19-fp-l.20 Saeid Nooshabad



Floating Point Number Distribution

•Which numbers can be represented?

-	8	-7	-6	-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	
					e tra	a ba	a huu	ud I		- hui	uul eu	i bi	. I .	• 4 ·				
								m							<u> </u>			
ĉ	4				22		3	3	201	NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN	3		22	222	222	2020	^с	
>	<				\times		\times	\times	ŝ	××	×		×	$\times \times$	××	××	××	
	000-1-				-1.000		-1.000	-1.000	1.100	1.000	1.000		1.000	1.001	1.100	1.101	1.111	

Using mantissa of 1.000 and positive exponents

and Sign bit

and negative exponents

in each of the intervals of exponentially increasing size can represent 2^s (s=3 here) numbers of uniform difference

But how do we represent 0? Next Lecture!

"And in Conclusion.."

- ^oNumber of digits allocated to significand and exponent, and choice of base, can affect both the number of different representable values and the range of values.
- [°]Finite precision means we have to cope with roundoff error (arithmetic with inexact values) and truncation error (large values overwhelming small ones).
- IEEE 754 Standard allows Single Precision (1 word) and Double Precision (2-word) representation of FP. Nos.