

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point

CSE 351 Winter 2018

Instructor:
Mark Wyse

Teaching Assistants:
Kevin Bi Parker, DeWilde, Emily Furst,
Sarah House, Waylon Huang, Vinny Palaniappan

<http://skcd.com/571/>

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Administrivia

- ❖ Lab 1 due Friday (1/19)
 - Submit `bits.c` and `pointer.c`
- ❖ Homework 2 out since 1/15, due 1/24
 - On Integers, Floating Point, and x86-64

2

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Unsigned Multiplication in C

Operands: w bits u v

True Product: $2w$ bits $u \cdot v$

Discard w bits: w bits $\text{UMult}_w(u, v)$

- ❖ Standard Multiplication Function
 - Ignores high order w bits
- ❖ Implements Modular Arithmetic
 - $\text{UMult}_w(u, v) = u \cdot v \bmod 2^w$

3

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Multiplication with shift and add

- ❖ Operation $u \ll k$ gives $u \cdot 2^k$
 - Both signed and unsigned

Operands: w bits u 2^k

True Product: $w + k$ bits $u \cdot 2^k$

Discard k bits: w bits $\text{UMult}_w(u, 2^k)$

- ❖ Examples:
 - $u \ll 3 == u * 8$
 - $u \ll 5 - u \ll 3 == u * 24$
 - Most machines shift and add faster than multiply
 - *Compiler generates this code automatically*

4

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Number Representation Revisited

- ❖ What can we represent in one word?
 - Signed and Unsigned Integers
 - Characters (ASCII)
 - Addresses
- ❖ How do we encode the following:
 - Real numbers (e.g. 3.14159)
 - Very large numbers (e.g. 6.02×10^{23})
 - Very small numbers (e.g. 6.626×10^{-34})
 - Special numbers (e.g. ∞ , NaN)

} Floating Point

5

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Topics

- ❖ **Fractional binary numbers**
- ❖ IEEE floating-point standard
- ❖ Floating-point operations and rounding
- ❖ Floating-point in C

- ❖ There are many more details that we won't cover
 - It's a 58-page standard...

6

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

Representation of Fractions

- “Binary Point,” like decimal point, signifies boundary between integer and fractional parts:

Example 6-bit representation: $xx.yyyy$

Diagram showing bits $x^1, x^0, y^1, y^2, y^3, y^4$ with weights $2^1, 2^0, 2^{-1}, 2^{-2}, 2^{-3}, 2^{-4}$.

- Example:** $10.1010_2 = 1 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 1 \times 2^{-2} = 2.625_{10}$
- Binary point numbers that match the 6-bit format above range from 0 (00.0000₂) to 3.9375 (11.1111₂)

7

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

Scientific Notation (Decimal)

Diagram of $6.02_{10} \times 10^{23}$ with labels: mantissa, decimal point, exponent, radix (base).

- Normalized form:** exactly one digit (non-zero) to left of decimal point
- Alternatives to representing 1/1,000,000,000
 - Normalized:** 1.0×10^{-9}
 - Not normalized:** $0.1 \times 10^{-8}, 10.0 \times 10^{-10}$

8

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

Scientific Notation (Binary)

Diagram of $1.01_2 \times 2^{-1}$ with labels: mantissa, binary point, exponent, radix (base).

- Computer arithmetic that supports this called **floating point** due to the “floating” of the binary point
 - Declare such variable in C as `float` (or `double`)

9

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

Scientific Notation Translation




- Convert from scientific notation to binary point to decimal
 - Perform the multiplication by shifting the decimal until the exponent disappears
 - Example:** $1.011_2 \times 2^4 = 10110_2 = 22_{10}$
 - Example:** $1.011_2 \times 2^{-2} = 0.01011_2 = 0.34375_{10}$
- Convert from binary point to **normalized** scientific notation
 - Distribute out exponents until binary point is to the right of a single digit
 - Example:** $1101.001_2 = 1.101001_2 \times 2^3$

10

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

Floating Point Topics

- Fractional binary numbers
- IEEE floating-point standard**
- Floating-point operations and rounding
- Floating-point in C

- There are many more details that we won't cover
 - It's a 58-page standard...

11

UNIVERSITY of WASHINGTON | L66: Floating Point | CSE351, Winter 2018

IEEE Floating Point

- IEEE 754
 - Established in 1985 as uniform standard for floating point arithmetic
 - Main idea: make numerically sensitive programs portable
 - Specifies two things: representation and result of floating operations
 - Now supported by all major CPUs
- Driven by numerical concerns
 - Scientists/numerical analysts** want them to be as **real** as possible
 - Engineers** want them to be **easy to implement** and **fast**
 - In the end:
 - Scientists mostly won out
 - Nice standards for rounding, overflow, underflow, but...
 - Hard to make fast in hardware
 - Float operations can be an order of magnitude slower than integer ops**

12

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Encoding

- Use normalized, base 2 scientific notation:
 - Value: $\pm 1 \times \text{Mantissa} \times 2^{\text{Exponent}}$
 - Bit Fields: $(-1)^S \times 1.M \times 2^{(E-\text{bias})}$
- Representation Scheme:
 - Sign bit (0 is positive, 1 is negative)
 - Mantissa (a.k.a. significand) is the fractional part of the number in normalized form and encoded in bit vector **M**
 - Exponent weights the value by a (possibly negative) power of 2 and encoded in the bit vector **E**

1 bit 8 bits 23 bits

13

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

The Exponent Field

- Use **biased notation**
 - Read exponent as unsigned, but with *bias* of $2^{w-1}-1 = 127$
 - Representable exponents roughly ½ positive and ½ negative
 - Exponent 0 (Exp = 0) is represented as E = 0b 0111 1111
- Why biased?
 - Makes floating point arithmetic easier
 - Makes somewhat compatible with two's complement
- Practice: To encode in biased notation, add the bias then encode in unsigned:
 - Exp = 1 → E = 0b
 - Exp = 127 → E = 0b
 - Exp = -63 → E = 0b

14

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

The Mantissa (Fraction) Field

1 bit 8 bits 23 bits

$$(-1)^S \times (1.M) \times 2^{(E-\text{bias})}$$

- Note the implicit 1 in front of the M bit vector
 - Example: 0b 0011 1111 1100 0000 0000 0000 0000 0000 is read as $1.1_2 = 1.5_{10}$, *not* $0.1_2 = 0.5_{10}$
 - Gives us an extra bit of *precision*
- Mantissa "limits"
 - Low values near M = 0b0...0 are close to 2^{Exp}
 - High values near M = 0b1...1 are close to $2^{\text{Exp}+1}$

15

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Peer Instruction Question

- What is the correct value encoded by the following floating point number?
 - 0b 0 10000000 110000000000000000000000

A. + 0.75
 B. + 1.5
 C. + 2.75
 D. + 3.5
 E. We're lost...

16

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Precision and Accuracy

- Precision** is a count of the number of bits in a computer word used to represent a value
 - Capacity for accuracy
- Accuracy** is a measure of the difference between the *actual value of a number* and its computer representation
 - High precision permits high accuracy but doesn't guarantee it. It is possible to have high precision but low accuracy.
 - Example: float pi = 3.14;
 - pi will be represented using all 24 bits of the mantissa (highly precise), but is only an approximation (not accurate)

17

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Need Greater Precision?

- Double Precision** (vs. Single Precision) in 64 bits
 - Diagram: 64-bit format with sign bit (S) at bit 63, exponent field (E) of 11 bits from bit 62 to bit 52, and mantissa field (M) of 52 bits from bit 51 to bit 32. The mantissa field is also shown as 32 bits from bit 31 to bit 0.
- C variable declared as double
 - Exponent bias is now $2^{10}-1 = 1023$
 - Advantages: greater precision (larger mantissa), greater range (larger exponent)
 - Disadvantages: more bits used, slower to manipulate

18

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Representing Very Small Numbers

- But wait... what happened to zero?
 - Using standard encoding $0x00000000 =$
 - Special case:** E and M all zeros = 0
 - Two zeros! But at least $0x00000000 = 0$ like integers
 - New numbers closest to 0:
 - $a = 1.0...0_2 \times 2^{-126} = 2^{-126}$
 - $b = 1.0...01_2 \times 2^{-126} = 2^{-126} + 2^{-149}$
 - Normalization and implicit 1 are to blame
 - Special case:** $E = 0, M \neq 0$ are **denormalized numbers**

19

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Other Special Cases

- $E = 0xFF, M = 0: \pm \infty$
 - e.g. division by 0
 - Still work in comparisons!
- $E = 0xFF, M \neq 0: \text{Not a Number (NaN)}$
 - e.g. square root of negative number, $0/0, \infty - \infty$
 - NaN propagates through computations
 - Value of M can be useful in debugging
- New largest value (besides ∞)?
 - $E = 0xFF$ has now been taken!
 - $E = 0xFE$ has largest: $1.1...1_2 \times 2^{127} = 2^{128} - 2^{104}$

20

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Summary

- Floating point approximates real numbers:
 - Diagram: S (1 bit), E (8 bits), M (23 bits)
 - Handles large numbers, small numbers, special numbers
 - Exponent in biased notation (bias = $2^{w-1}-1$)
 - Outside of representable exponents is *overflow* and *underflow*
 - Mantissa approximates fractional portion of binary point
 - Implicit leading 1 (normalized) except in special cases
 - Exceeding length causes *rounding*

Exponent	Mantissa	Meaning
0x00	0	± 0
0x00	non-zero	\pm denorm num
0x01 - 0xFE	anything	\pm norm num
0xFF	0	$\pm \infty$
0xFF	non-zero	NaN

21

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating point topics

- Fractional binary numbers
- IEEE floating-point standard
- Floating-point operations and rounding**
- Floating-point in C
- There are many more details that we won't cover
 - It's a 58-page standard...

22

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Distribution of Values

- What ranges are NOT representable?
 - Between largest norm and infinity **Overflow**
 - Between zero and smallest denorm **Underflow**
 - Between norm numbers? **Rounding**
- Given a FP number, what's the bit pattern of the next largest representable number?
 - What is this "step" when $Exp = 0$?
 - What is this "step" when $Exp = 100$?
- Distribution of values is denser toward zero
 - Diagram: Number line from -15 to 15. Blue diamonds (Denormalized) are clustered near 0. Purple triangles (Normalized) are spread out. Yellow squares (Infinity) are at the ends.

23

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Operations: Basic Idea

Value = $(-1)^s \times \text{Mantissa} \times 2^{\text{Exponent}}$

Diagram: S (1 bit), E (8 bits), M (23 bits)

- $x +_f y = \text{Round}(x + y)$
- $x *_f y = \text{Round}(x * y)$
- Basic idea for floating point operations:
 - First, **compute the exact result**
 - Then **round** the result to make it fit into desired precision:
 - Possibly over/underflow if exponent outside of range
 - Possibly drop least-significant bits of mantissa to fit into M bit vector

24

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Addition Line up the binary points!

$(-1)^{S1} \times \text{Man1} \times 2^{\text{Exp1}} + (-1)^{S2} \times \text{Man2} \times 2^{\text{Exp2}}$

- Assume $\text{Exp1} > \text{Exp2}$

$$\begin{array}{r} 1.010 \times 2^2 \\ + 1.000 \times 2^{-1} + 0.0010 \times 2^2 \\ \hline \end{array}$$
- Exact Result: $(-1)^S \times \text{Man} \times 2^{\text{Exp}}$
 - Sign S , mantissa Man :
 - Result of signed align & add
 - Exponent E : $E1$
- Adjustments:
 - If $\text{Man} \geq 2$, shift Man right, increment Exp
 - If $\text{Man} < 1$, shift Man left k positions, decrement Exp by k
 - Over/underflow if Exp out of range
 - Round Man to fit mantissa precision

25

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Multiplication

$(-1)^{S1} \times \text{Man1} \times 2^{\text{Exp1}} \times (-1)^{S2} \times \text{Man2} \times 2^{\text{Exp2}}$

- Exact Result: $(-1)^S \times \text{Man} \times 2^E$
 - Sign S : $S1 \wedge S2$
 - Mantissa Man : $\text{Man1} \times \text{Man2}$
 - Exponent Exp : $\text{Exp1} + \text{Exp2}$
- Adjustments:
 - If $\text{Man} \geq 2$, shift Man right, increment Exp
 - Over/underflow if Exp out of range
 - Round Man to fit mantissa precision

26

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Mathematical Properties of FP Operations

- Exponent overflow yields $+\infty$ or $-\infty$
- Floats with value $+\infty$, $-\infty$, and NaN can be used in operations
 - Result usually still $+\infty$, $-\infty$, or NaN; but not always intuitive
- Floating point operations do not work like real math, due to **rounding**
 - Not associative: $(3.14 + 1e100) - 1e100 \neq 3.14 + (1e100 - 1e100)$

$$\begin{array}{r} 0 \\ 3.14 \end{array}$$
 - Not distributive: $100 * (0.1 + 0.2) \neq 100 * 0.1 + 100 * 0.2$

$$\begin{array}{r} 30.00000000000003553 \\ 30 \end{array}$$
 - Not cumulative
 - Repeatedly adding a very small number to a large one may do nothing

27

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating point topics

- Fractional binary numbers
- IEEE floating-point standard
- Floating-point operations and rounding
- Floating-point in C**
- There are many more details that we won't cover
 - It's a 58-page standard...

28

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point in C !!!

- C offers two (well, 3) levels of precision

float	1.0f	single precision (32-bit)
double	1.0	double precision (64-bit)
long double	1.0L	("double double" or quadruple) precision (64-128 bits)
- `#include <math.h>` to get INFINITY and NAN constants
- Equality (`==`) comparisons between floating point numbers are tricky, and often return unexpected results, so just avoid them!

29

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Conversions in C !!!

- Casting between int, float, and double changes the bit representation**
 - `int` \rightarrow `float`
 - May be rounded (not enough bits in mantissa: 23)
 - Overflow impossible
 - `int` or `float` \rightarrow `double`
 - Exact conversion (all 32-bit ints representable)
 - `long` \rightarrow `double`
 - Depends on word size (32-bit is exact, 64-bit may be rounded)
 - `double` or `float` \rightarrow `int`
 - Truncates fractional part (rounded toward zero)
 - "Not defined" when out of range or NaN: generally sets to `Tmin` (even if the value is a very big positive)

30

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Number Representation Really Matters

- ❖ **1991:** Patriot missile targeting error
 - clock skew due to conversion from integer to floating point
- ❖ **1996:** Ariane 5 rocket exploded (\$1 billion)
 - overflow converting 64-bit floating point to 16-bit integer
- ❖ **2000:** Y2K problem
 - limited (decimal) representation: overflow, wrap-around
- ❖ **2038:** Unix epoch rollover
 - Unix epoch = seconds since 12am, January 1, 1970
 - signed 32-bit integer representation rolls over to TMin in 2038
- ❖ **Other related bugs:**
 - 1982: Vancouver Stock Exchange 10% error in less than 2 years
 - 1994: Intel Pentium FDIV (floating point division) HW bug (\$475 million)
 - 1997: USS Yorktown "smart" warship stranded: divide by zero
 - 1998: Mars Climate Orbiter crashed: unit mismatch (\$193 million)

31

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Summary

- ❖ Floats also suffer from the fixed number of bits available to represent them
 - Can get overflow/underflow
 - "Gaps" produced in representable numbers means we can lose precision, unlike `ints`
 - Some "simple fractions" have no exact representation (e.g. 0.2)
 - "Every operation gets a slightly wrong result"
- ❖ Floating point arithmetic not associative or distributive
 - Mathematically equivalent ways of writing an expression may compute different results
- ❖ **Never** test floating point values for equality!
- ❖ **Careful** when converting between `ints` and `floats`!

32

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point Summary

- ❖ Converting between integral and floating point data types *does* change the bits
 - Floating point rounding is a HUGE issue!
 - Limited mantissa bits cause inaccurate representations
 - Floating point arithmetic is NOT associative or distributive

33

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Denorm Numbers

This is extra (non-testable) material

- ❖ Denormalized numbers
 - No leading 1
 - Uses implicit exponent of -126 even though `E = 0x00`
- ❖ Denormalized numbers close the gap between zero and the smallest normalized number
 - Smallest norm: $\pm 1.0...0_{\text{two}} \times 2^{-126} = \pm 2^{-126}$
 - Smallest denorm: $\pm 0.0...01_{\text{two}} \times 2^{-126} = \pm 2^{-149}$
 - There is still a gap between zero and the smallest denormalized number

So much closer to 0

34

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Floating Point and the Programmer

```

#include <stdio.h>

int main(int argc, char* argv[]){
    float f1 = 1.0;
    float f2 = 0.0;
    int i;
    for (i = 0; i < 10; i++)
        f2 += 1.0/10.0;

    printf("0x%08x 0x%08x\n", *(int*)&f1, *(int*)&f2);
    printf("f1 = %10.9f\n", f1);
    printf("f2 = %10.9f\n", f2);

    f1 = 1E30;
    f2 = 1E-30;
    float f3 = f1 + f2;
    printf("f1 == f3? %s\n", f1 == f3 ? "yes" : "no" );

    return 0;
}

```

```

$ ./a.out
0x3f800000 0x3f800001
f1 = 1.000000000
f2 = 1.000000119
f1 == f3? yes

```

35

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

BONUS SLIDES

An example that applies the IEEE Floating Point concepts to a smaller (8-bit) representation scheme.

36

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Tiny Floating Point Example

- ❖ 8-bit Floating Point Representation
 - The sign bit is in the most significant bit (MSB)
 - The next four bits are the exponent, with a bias of $2^{4-1}-1 = 7$
 - The last three bits are the mantissa
- ❖ Same general form as IEEE Format
 - Normalized binary scientific point notation
 - Similar special cases for 0, denormalized numbers, NaN, ∞

37

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Dynamic Range (Positive Only)

	S	E	M	Exp	Value	
Denormalized numbers	0	0000	000	-6	0	
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	closest to zero
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	...					
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	
Normalized numbers	0	0000	111	-6	$7/8 * 1/64 = 7/512$	largest denorm
	0	0001	000	-6	$8/8 * 1/64 = 8/512$	smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$	
	...					
	0	0110	110	-1	$14/8 * 1/2 = 14/16$	
	0	0110	111	-1	$15/8 * 1/2 = 15/16$	closest to 1 below
	0	0111	000	0	$8/8 * 1 = 1$	
	0	0111	001	0	$9/8 * 1 = 9/8$	closest to 1 above
	0	0111	010	0	$10/8 * 1 = 10/8$	
	...					
0	1110	110	7	$14/8 * 128 = 224$		
0	1110	111	7	$15/8 * 128 = 240$	largest norm	
0	1111	000	n/a	inf		

38

UNIVERSITY of WASHINGTON L06: Floating Point CSE351, Winter 2018

Special Properties of Encoding

- ❖ Floating point zero (0^+) exactly the same bits as integer zero
 - All bits = 0
- ❖ Can (Almost) Use Unsigned Integer Comparison
 - Must first compare sign bits
 - Must consider $0^- = 0^+ = 0$
 - NaNs problematic
 - Will be greater than any other values
 - What should comparison yield?
 - Otherwise OK
 - Denorm vs. normalized
 - Normalized vs. infinity

39