

18-600 Foundations of Computer Systems

Lecture 4: “Floating Point”

John Shen & Zhiyi Yu

September 12, 2016

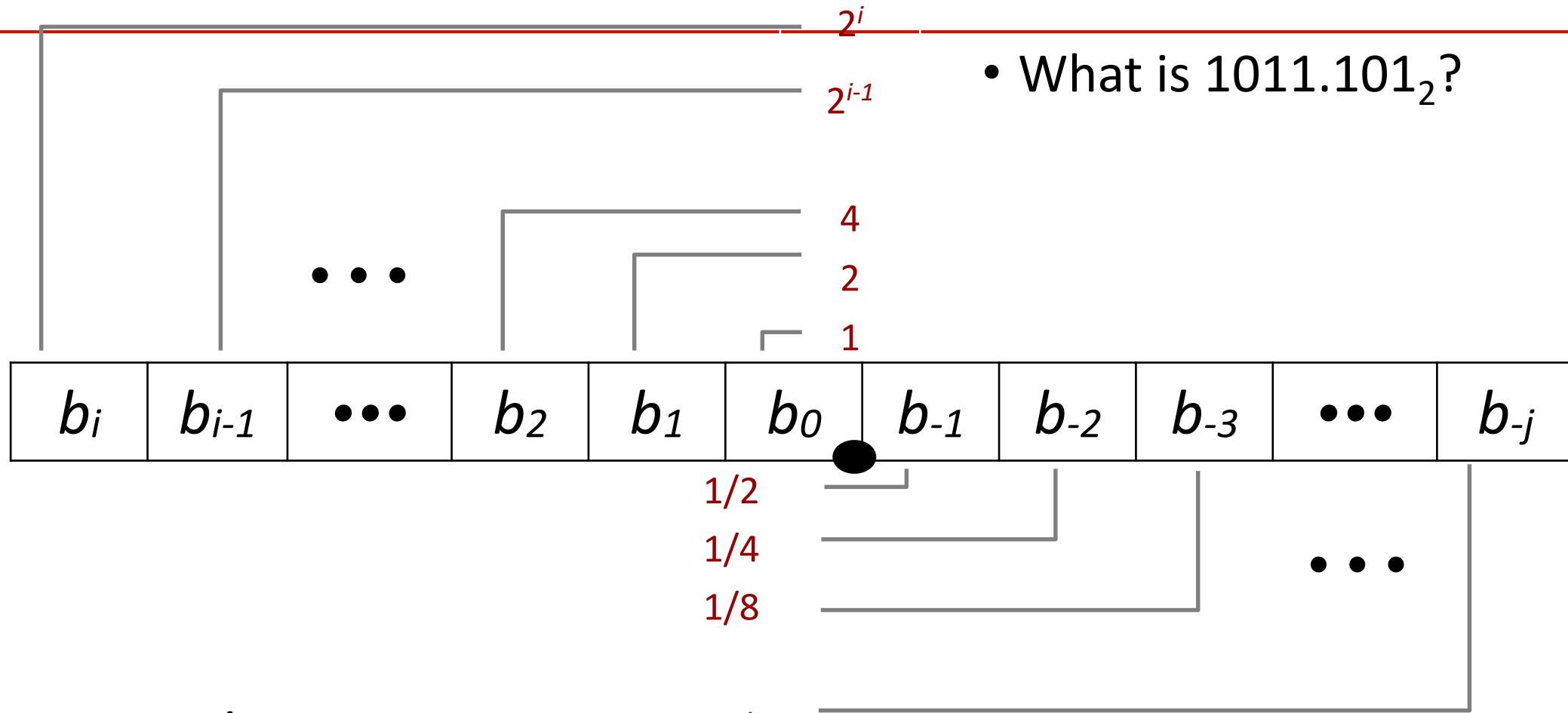
- Required Reading Assignment:
 - Chapter 2 of CS:APP (3rd edition) by Randy Bryant & Dave O'Hallaron
- Assignments for This Week:
 - ❖ Lab 1



Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

Fractional Binary Numbers



- Representation
 - Bits to right of “binary point” represent fractional powers of 2
 - Represents rational number:

$$\sum_{k=-j}^i b_k \times 2^k$$

Fractional Binary Numbers: Examples



Value

5 3/4

2 7/8

1 7/16

Representation

101.11_2

10.111_2

1.0111_2



Observations

- Divide by 2 by shifting right (unsigned)
- Multiply by 2 by shifting left
- Numbers of form $0.111111\dots_2$ are just below 1.0
 - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
 - Use notation $1.0 - \epsilon$

Representable Numbers

- Limitation #1
 - Can only exactly represent numbers of the form $x/2^k$
 - Other rational numbers have repeating bit representations
 - Value Representation
 - $1/3$ 0.0101010101 [01]...₂
 - $1/5$ 0.001100110011 [0011]...₂
 - $1/10$ 0.0001100110011 [0011]...₂
- Limitation #2
 - Just one setting of binary point within the w bits
 - Limited range of numbers (very small values? very large?)

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

IEEE Floating Point

- IEEE Standard 754
 - Established in 1985 as uniform standard for floating point arithmetic
 - Before that, many idiosyncratic formats
 - Supported by all major CPUs
- Driven by numerical concerns
 - Nice standards for rounding, overflow, underflow
 - Hard to make fast in hardware
 - Numerical analysts predominated over hardware designers in defining standard

Floating Point Representation

- Numerical Form:

$$(-1)^s M 2^E$$

- **Sign bit s** determines whether number is negative or positive
- **Significand M** (mantissa) normally a fractional value in range [1.0,2.0).
- **Exponent E** weights value by power of two

- Encoding

- MSB s is sign bit s
- exp field encodes E (but is not equal to E)
- frac field encodes M (but is not equal to M)



Precision options

- Single precision: 32 bits



- Double precision: 64 bits



- Extended precision: 80 bits (Intel only)



"Normalized" Values

$$v = (-1)^s M 2^E$$

- When: $\text{exp} \neq 000\ldots0$ and $\text{exp} \neq 111\ldots1$
- Exponent coded as a ***biased*** value: $E = \text{Exp} - \text{Bias}$
 - Exp : unsigned value of exp field
 - $\text{Bias} = 2^{k-1} - 1$, where k is number of exponent bits
 - Single precision: 127 ($\text{Exp}: 1\ldots254$, $E: -126\ldots127$)
 - Double precision: 1023 ($\text{Exp}: 1\ldots2046$, $E: -1022\ldots1023$)
- Significand coded with implied leading 1: $M = 1.\text{xxx}\ldots\text{x}_2$
 - $\text{xxx}\ldots\text{x}$: bits of frac field
 - Minimum when $\text{frac}=000\ldots0$ ($M = 1.0$)
 - Maximum when $\text{frac}=111\ldots1$ ($M = 2.0 - \epsilon$)
 - Get extra leading bit for "free"

Normalized Encoding Example

$$v = (-1)^s M 2^E$$

$$E = Exp - Bias$$

- Value: float F = 18600.0;
 - $18600_{10} = 100100010101000_2 = 1.00100010101_2 \times 2^{14}$

- Significand

$M = 1.\underline{00100010101}_2$
 $\text{frac} = \underline{00100010101}000000000000_2 \quad (23 \text{ bits})$

- Exponent

$E = 14$
 $Bias = 127$
 $Exp = 141 = 10001101_2 \quad (8 \text{ bits})$

- Result:

0	10001101	001000010101000000000000
s	exp	frac

Denormalized Values

$$v = (-1)^s M 2^E$$

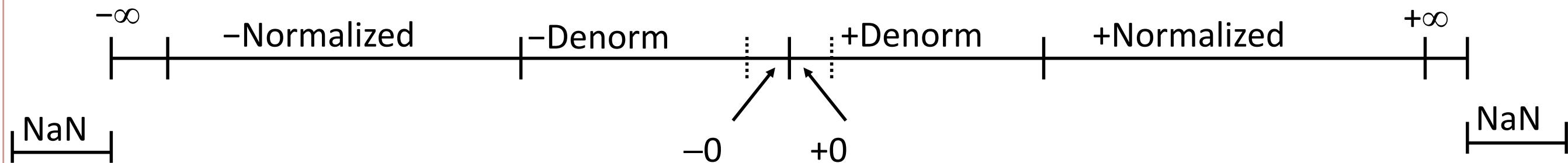
$$E = 1 - \text{Bias}$$

- Condition: $\text{exp} = 000\dots0$
- Exponent value: $E = 1 - \text{Bias}$ (instead of $E = 0 - \text{Bias}$)
- Significand coded with implied leading 0: $M = 0.\text{xxx}\dots\text{x}_2$
 - $\text{xxx}\dots\text{x}$: bits of **frac**
- Cases
 - **exp** = 000...0, **frac** = 000...0
 - Represents zero value
 - Note distinct values: +0 and -0 (why?)
 - **exp** = 000...0, **frac** ≠ 000...0
 - Numbers closest to 0.0
 - Equispaced

Special Values

- Condition: **exp = 111...1**
- Case: **exp = 111...1, frac = 000...0**
 - Represents value ∞ (infinity)
 - Operation that overflows
 - Both positive and negative
 - E.g., $1.0/0.0 = -1.0/-0.0 = +\infty$, $1.0/-0.0 = -\infty$
- Case: **exp = 111...1, frac \neq 000...0**
 - Not-a-Number (NaN)
 - Represents case when no numeric value can be determined
 - E.g., $\sqrt{-1}$, $\infty - \infty$, $\infty \times 0$

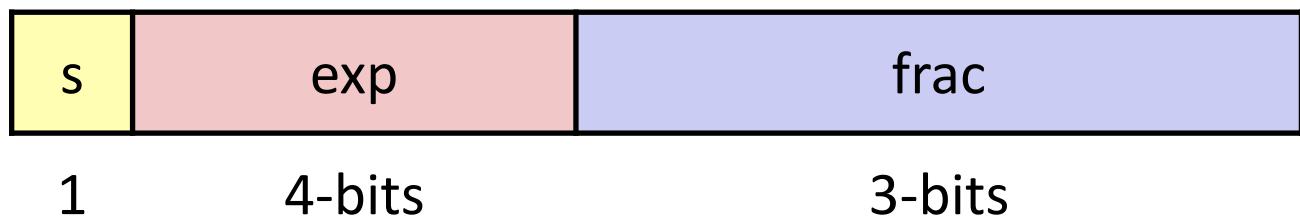
Visualization: Floating Point Encodings



Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

Tiny Floating Point Example



- 8-bit Floating Point Representation
 - the sign bit is in the most significant bit
 - the next four bits are the exponent, with a bias of 7
 - the last three bits are the **frac**
- Same general form as IEEE Format
 - normalized, denormalized
 - representation of 0, NaN, infinity

Dynamic Range (Positive Only)

	s	exp	frac	E	Value	
Denormalized numbers	0	0000	000	-6	0	
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	...					closest to zero
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	
	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	$14/8 * 1/2 = 14/16$	
Normalized numbers	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	

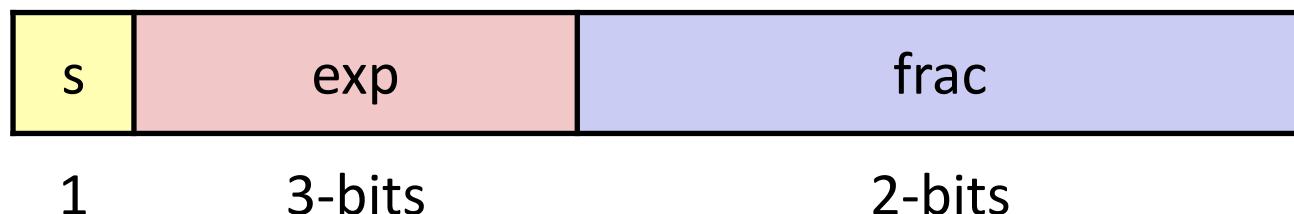
$$v = (-1)^s M 2^E$$

$$n: E = \text{Exp} - \text{Bias}$$

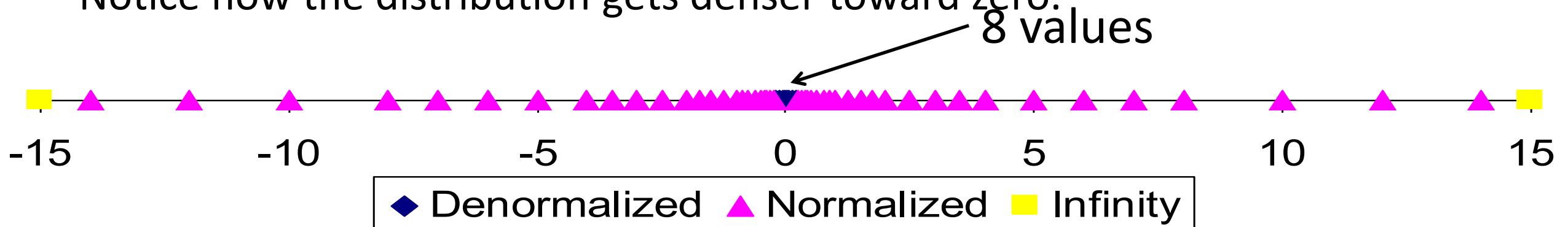
$$d: E = 1 - \text{Bias}$$

Distribution of Values

- 6-bit IEEE-like format
 - $e = 3$ exponent bits
 - $f = 2$ fraction bits
 - Bias is $2^{3-1}-1 = 3$



- Notice how the distribution gets denser toward zero.

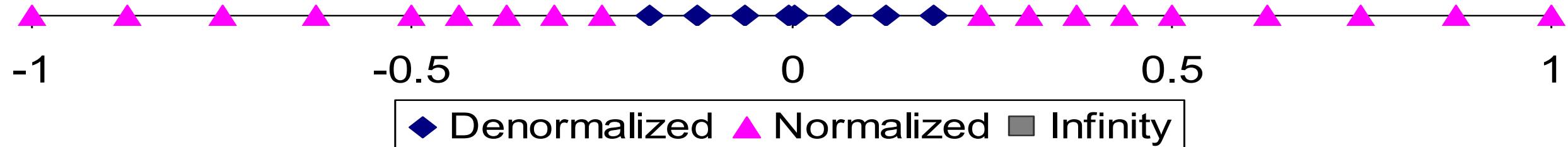


Distribution of Values (close-up view)

- 6-bit IEEE-like format
 - $e = 3$ exponent bits
 - $f = 2$ fraction bits
 - Bias is 3



1 3-bits 2-bits



Special Properties of the IEEE Encoding

- FP Zero Same as Integer Zero
 - All bits = 0
- Can (Almost) Use Unsigned Integer Comparison
 - Must first compare sign bits
 - Must consider $-0 = 0$
 - NaNs problematic
 - Will be greater than any other values
 - What should comparison yield?
 - Otherwise OK
 - Denorm vs. normalized
 - Normalized vs. infinity

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- **Rounding, addition, multiplication**
- Floating point in C
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

Floating Point Operations: Basic Idea

- $x +_f y = \text{Round}(x + y)$
- $x \times_f y = \text{Round}(x \times y)$
- Basic idea
 - First compute exact result
 - Make it fit into desired precision
 - Possibly overflow if exponent too large
 - Possibly round to fit into `frac`

Rounding

- Rounding Modes (illustrate with \$ rounding)

•	\$1.40	\$1.60	\$1.50	\$2.50	-\$1.50
• Towards zero	\$1	\$1	\$1	\$2	-\$1
• Round down ($-\infty$)	\$1	\$1	\$1	\$2	-\$2
• Round up ($+\infty$)	\$2	\$2	\$2	\$3	-\$1
• Nearest Even (default)	\$1	\$2	\$2	\$2	-\$2

Closer Look at Round-To-Even

- Default Rounding Mode
 - Hard to get any other kind without dropping into assembly
 - All others are statistically biased
 - Sum of set of positive numbers will consistently be over- or under- estimated
- Applying to Other Decimal Places / Bit Positions
 - When exactly halfway between two possible values
 - Round so that least significant digit is even
 - E.g., round to nearest hundredth

7.8949999	7.89	(Less than half way)
7.8950001	7.90	(Greater than half way)
7.8950000	7.90	(Half way—round up)
7.8850000	7.88	(Half way—round down)

Rounding Binary Numbers

- Binary Fractional Numbers
 - “Even” when least significant bit is 0
 - “Half way” when bits to right of rounding position = $100\dots_2$

- Examples
 - Round to nearest $1/4$ (2 bits right of binary point)

Value	Binary	Rounded	Action	Rounded Value
$2\frac{3}{32}$	$10.00\textcolor{red}{011}_2$	10.00_2	($<1/2$ —down)	2
$2\frac{3}{16}$	$10.00\textcolor{red}{110}_2$	10.01_2	($>1/2$ —up)	$2\frac{1}{4}$
$2\frac{7}{8}$	$10.11\textcolor{red}{100}_2$	11.00_2	($1/2$ —up)	3
$2\frac{5}{8}$	$10.10\textcolor{red}{100}_2$	10.10_2	($1/2$ —down)	$2\frac{1}{2}$

FP Multiplication

- $(-1)^{s_1} M_1 2^{E_1} \times (-1)^{s_2} M_2 2^{E_2}$
- Exact Result: $(-1)^s M 2^E$
 - Sign s : $s_1 \wedge s_2$
 - Significand M : $M_1 \times M_2$
 - Exponent E : $E_1 + E_2$
- Fixing
 - If $M \geq 2$, shift M right, increment E
 - If E out of range, overflow
 - Round M to fit `frac` precision
- Implementation
 - Biggest chore is multiplying significands

Floating Point Addition

- $(-1)^{s1} M1 2^{E1} + (-1)^{s2} M2 2^{E2}$

- Assume $E1 > E2$

- Exact Result: $(-1)^s M 2^E$

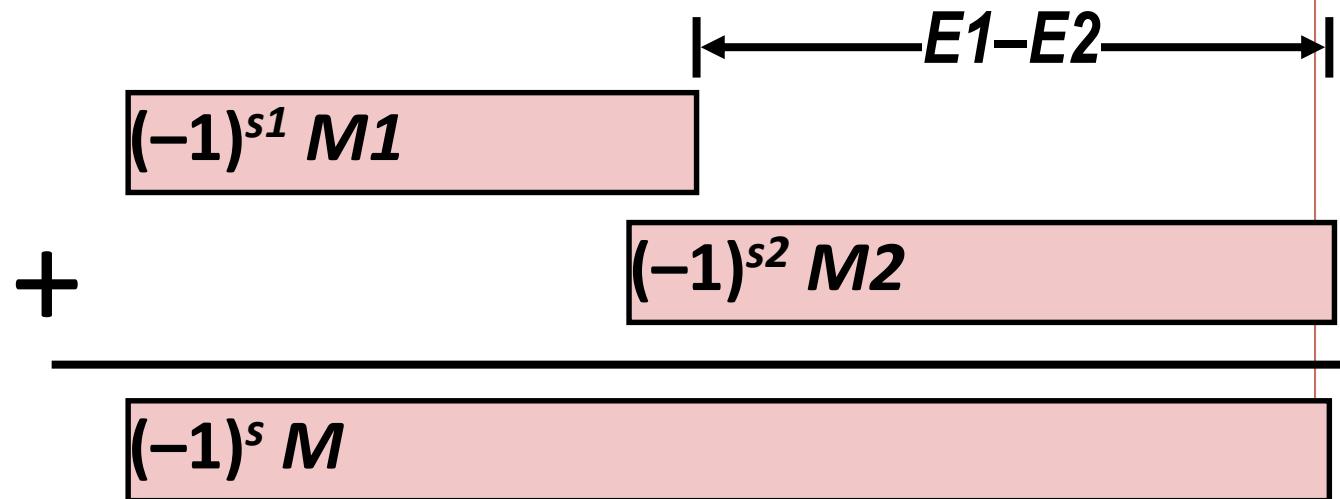
- Sign s , significand M :
 - Result of signed align & add

- Exponent E : $E1$

- Fixing

- If $M \geq 2$, shift M right, increment E
- If $M < 1$, shift M left k positions, decrement E by k
- Overflow if E out of range
- Round M to fit **frac** precision

Get binary points lined up



Mathematical Properties of FP Add

- Commutative? **Yes**
- Associative?
 - Overflow and inexactness of rounding
 - $(3.14 + 1e10) - 1e10 = 0, 3.14 + (1e10 - 1e10) = 3.14$**No**
- 0 is additive identity? **Yes**
- Every element has additive inverse?
 - Yes, except for infinities & NaNs**Almost**
- Monotonicity
 - $a \geq b \Rightarrow a+c \geq b+c?$
 - Except for infinities & NaNs**Almost**

Mathematical Properties of FP Mult

- Multiplication Commutative? **Yes**
- Multiplication is Associative?
 - Possibility of overflow, inexactness of rounding
 - Ex: $(1e20 * 1e20) * 1e-20 = \inf$, $1e20 * (1e20 * 1e-20) = 1e20$**No**
- 1 is multiplicative identity? **Yes**
- Multiplication distributes over addition?
 - Possibility of overflow, inexactness of rounding
 - $1e20 * (1e20 - 1e20) = 0.0$, $1e20 * 1e20 - 1e20 * 1e20 = \text{NaN}$**No**
- Monotonicity
 - $a \geq b \ \& \ c \geq 0 \Rightarrow a * c \geq b * c?$
 - Except for infinities & NaNs**Almost**

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- **Floating point in C**
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

Floating Point in C

- C Guarantees Two Levels
 - **float** single precision
 - **double** double precision
- Conversions/Casting
 - Casting between **int**, **float**, and **double** changes bit representation
 - **double/float → int**
 - Truncates fractional part
 - Like rounding toward zero
 - Not defined when out of range or NaN: Generally sets to TMin
 - **int → double**
 - Exact conversion, as long as **int** has \leq 53 bit word size
 - **int → float**
 - Will round according to rounding mode

Floating Point Puzzles

- For each of the following C expressions, either:

- Argue that it is true for all argument values
- Explain why not true

```
int x = ...;
float f = ...;
double d = ...;
```

Assume neither
 d nor f is NaN

- $x == (int)(float) x$ X
- $x == (int)(double) x$ ✓
- $f == (float)(double) f$ ✓
- $d == (double)(float) d$ X
- $f == -(-f);$ ✓
- $2/3 == 2/3.0$ X
- $d < 0.0 \Rightarrow ((d*2) < 0.0)$ ✓
- $d > f \Rightarrow -f > -d$ ✓
- $d * d >= 0.0$ ✓
- $(d+f)-d == f$ X

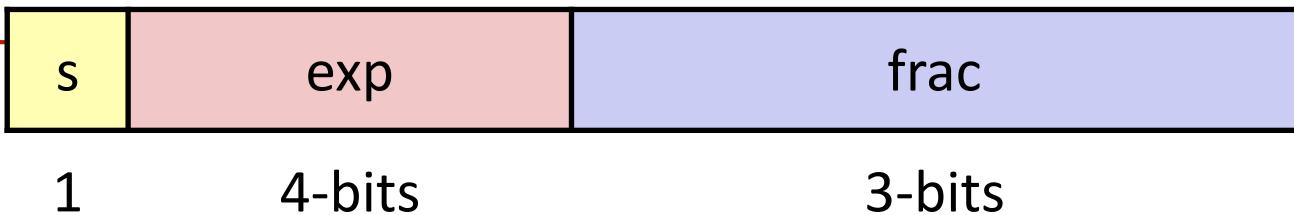
Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- **Creating floating point number**
- Summary of floating point number
- Quick introduction of assembly language

Creating Floating Point Number

- Steps

- Normalize to have leading 1
- Round to fit within fraction
- Postnormalize to deal with effects of rounding



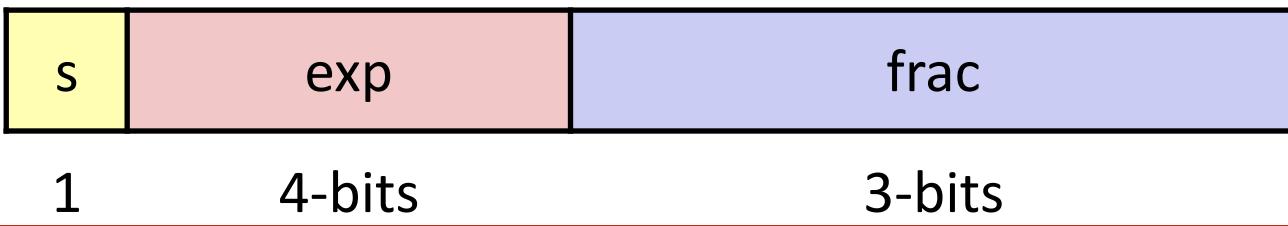
- Case Study

- Convert 8-bit unsigned numbers to tiny floating point format

Example Numbers

128	10000000
15	00001101
17	00010001
19	00010011
138	10001010
63	00111111

Normalize



- Requirement
 - Set binary point so that numbers of form 1.xxxxx
 - Adjust all to have leading one
 - Decrement exponent as shift left

Value	Binary	Fraction	Exponent
128	10000000	1.0000000	7
15	00001101	1.1010000	3
17	00010001	1.0001000	4
19	00010011	1.0011000	4
138	10001010	1.0001010	7
63	00111111	1.1111100	5

Rounding

Guard bit: LSB of result

Round bit: 1st bit removed

1 . BBGRXXXX

Sticky bit: OR of remaining bits

- Round up conditions
 - Round = 1, Sticky = 1 $\rightarrow > 0.5$
 - Guard = 1, Round = 1, Sticky = 0 \rightarrow Round to even

<i>Value</i>	<i>Fraction</i>	<i>GRS</i>	<i>Incr?</i>	<i>Rounded</i>
128	1.0000000	000	N	1.000
15	1.1010000	100	N	1.101
17	1.0001000	010	N	1.000
19	1.0011000	110	Y	1.010
138	1.0001010	011	Y	1.001
63	1.1111100	111	Y	10.000

Postnormalize

- Issue
 - Rounding may have caused overflow
 - Handle by shifting right once & incrementing exponent

Value	Rounded	Exp	Adjusted	Result
128	1.000	7		128
15	1.101	3		15
17	1.000	4		16
19	1.010	4		20
138	1.001	7		134
63	10.000	5	1.000/6	64

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Creating floating point number
- **Summary of floating point number**
- Quick introduction of assembly language

Summary

- IEEE Floating Point has clear mathematical properties
- Represents numbers of form $M \times 2^E$
- One can reason about operations independent of implementation
 - As if computed with perfect precision and then rounded
- Not the same as real arithmetic
 - Violates associativity/distributivity
 - Makes life difficult for compilers & serious numerical applications programmers

Today: Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Creating floating point number
- Summary of floating point number
- Quick introduction of assembly language

Compiling C Into Assembly

C Code (sum.c)

```
long plus(long x, long y);  
  
void sumstore(long x, long y,  
              long *dest)  
{  
    long t = plus(x, y);  
    *dest = t;  
}
```

Generated x86-64 Assembly

```
sumstore:  
    pushq   %rbx  
    movq   %rdx, %rbx  
    call    plus  
    movq   %rax, (%rbx)  
    popq   %rbx  
    ret
```

Obtain (on shark machine) with command

```
gcc -Og -S sum.c
```

Produces file sum.s

Machine Instruction Example

```
*dest = t;
```

- C Code
 - Store value **t** where designated by **dest**

```
movq %rax, (%rbx)
```

- Assembly
 - Move 8-byte value to memory
 - Quad words in x86-64 parlance
 - Operands:
 - t**:Register **%rax**
 - dest**: Register **%rbx**
 - *dest**: Memory **M[%rbx]**

```
0x40059e: 48 89 03
```

- Object Code
 - 3-byte instruction
 - Stored at address **0x40059e**

Assembly Characteristics: Operations

- Perform arithmetic function on register or memory data
- Transfer data between memory and register
 - Load data from memory into register
 - Store register data into memory
- Transfer control
 - Unconditional jumps to/from procedures
 - Conditional branches

Some Arithmetic Operations

- Two Operand Instructions:

Format	Computation	
addq	<i>Src, Dest</i> Dest = Dest + Src	
subq	<i>Src, Dest</i> Dest = Dest – Src	
imulq	<i>Src, Dest</i> Dest = Dest * Src	
salq	<i>Src, Dest</i> Dest = Dest << Src	<i>Also called shlq</i>
sarq	<i>Src, Dest</i> Dest = Dest >> Src	<i>Arithmetic</i>
shrq	<i>Src, Dest</i> Dest = Dest >> Src	<i>Logical</i>
xorq	<i>Src, Dest</i> Dest = Dest ^ Src	
andq	<i>Src, Dest</i> Dest = Dest & Src	
orq	<i>Src, Dest</i> Dest = Dest Src	

- One Operand Instructions

incq	<i>Dest</i>	Dest = Dest + 1
decq	<i>Dest</i>	Dest = Dest – 1
negq	<i>Dest</i>	Dest = – Dest
notq	<i>Dest</i>	Dest = ~Dest

x86-64 Integer Registers

%rax	%eax
%rbx	%ebx
%rcx	%ecx
%rdx	%edx
%rsi	%esi
%rdi	%edi
%rsp	%esp
%rbp	%ebp

%r8	%r8d
%r9	%r9d
%r10	%r10d
%r11	%r11d
%r12	%r12d
%r13	%r13d
%r14	%r14d
%r15	%r15d

Simple Memory Addressing Modes

- Normal (R) Mem[Reg[R]]
 - Register R specifies memory address
 - Pointer dereferencing in C

movq (%rcx), %rax

- Displacement D(R) Mem[Reg[R]+D]
 - Register R specifies start of memory region
 - Constant displacement D specifies offset

movq 8(%rbp), %rdx

18-600 Foundations of Computer Systems

Lecture 5: "Machine Programs I: (Basics)"

John Shen & Zhiyi Yu

September 14, 2016

Next Time ...



Electrical & Computer
ENGINEERING