Floating Point

(with contributions from Dr. Bin Ren, William & Mary Computer Science)

Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Fractional Binary Numbers



- 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

What is 1011.101₂?

1/2 + 1/8 = 5/8 11 5/8 or 11.625

What is 123.45 in binary?

- **123 = 1111011**
- to get the .45, use repeated multiplication by 2
 - if product < 1, bit is 0</p>
 - if product >= 1, bit is 1 and subtract 1 from product
 - .45 * 2 =

.9 * 2 =	0
(1.8-1) * 2 =	1
(1.6-1) * 2 =	1
(1.2 – 1) * 2 =	1
.4	0

1111011.01110

Fractional Binary Numbers: Examples

Value	Representation
5 3/4	101.112
2 7/8	10.1112
1 7/16	1.01112

Observations

- Divide by 2 by shifting right (unsigned)
 - compare 101.11₂ with 10.111₂
- Multiply by 2 by shifting left
 - compare 101.11₂ with 1011.1₂
- Numbers of form 0.111111...2 are just below 1.0
 - 1/2 + 1/4 + 1/8 + ... + 1/2ⁱ + ... → 1.0
 - notation sometimes seen: 1.0 ε

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?)

Representable Numbers

Numbers 0.111...11 base 2 represent numbers just below 1

0.111111 base 2 = 63/64

Only finite-length encodings

- 1/3 and 5/7 cannot be represented exactly
- Fractional binary notation can only represent numbers that can be written x * 2^y (i.e. 63/64 = 63*2⁻⁶)
 - Otherwise, approximated
 - Increasing accuracy = lengthening the binary representation but still have finite space

Practice

Fractional value of the following binary values:

- **.**01 =
- **.**010 =
- **1.00110 =**
- **11.001101 =**

123.45 base 10

- Binary value =
- FYI also equals:
 - 1.2345 x 10² in normalized form
 - 12345 x 10⁻² using significand/mantissa/coefficient and exponent

Floating Point

Background: Fractional binary numbers

IEEE floating point standard: Definition

- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

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
 - Intel-based PCs
 - Apple
 - Unix/Linux

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 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)

s	ехр	frac

Precision options

Single precision: 32 bits

S	ехр	frac
1	8-hits	23-hits

Double precision: 64 bits



Extended precision: 80 bits (Intel only)



Normalized Values

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

When: exp ≠ 000...0 and exp ≠ 111...1

Exponent coded as a *biased* value: E = Exp – Bias

- Exp: unsigned value of exp field
- Bias = 2^{k-1} 1, where k is number of exponent bits
 - Single precision: 127 (Exp: 1...254, E: -126...127)
 - Double precision: 1023 (Exp: 1...2046, E: -1022...1023)

Significand coded with implied leading 1: M = 1.xxx...x₂

- xxx...x: bits of frac field
- Minimum when frac = 000...0 (M = 1.0)
- Maximum when frac = 111...1 (M = 2.0ε)
- Get extra leading bit for "free"

Bias Notes

Biasing is done because exponents have to be signed values in order to be able to represent both tiny and huge values, but two's complement, the usual representation for signed values, would make comparison harder.

- To solve this problem the exponent is biased to put it within an unsigned range suitable for comparison.
- By arranging the fields so that the sign bit is in the most significant bit position, the biased exponent in the middle, then the mantissa in the least significant bits, the resulting value will be ordered properly, whether it's interpreted as a floating point or integer value. This allows high speed comparisons of floating point numbers using fixed point hardware.

When interpreting the floating-point number, the bias is subtracted to retrieve the actual exponent.

Significand Notes

- Represents the fraction, or precision bits of the number.
- It is composed of an implicit (i.e., hidden) leading bit and the fraction bits.
- In order to maximize the quantity of representable numbers, floating-point numbers are typically stored in *normalized* form.
 - This basically puts the radix point after the first non-zero digit
 - Nice optimization available in base two, since the only possible non-zero digit is 1.
 - Thus, we can just assume a leading digit of 1, and don't need to represent it explicitly.
 - As a result, the mantissa/significand has effectively 24 bits of resolution, by way of 23 fraction bits.

Normalized Encoding Example

 $v = (-1)^{s} M 2^{E}$ E = Exp - Bias

Value: float F = 15213.0;

15213₁₀ = 11101101101101₂ = 1.1101101101101₂ x 2¹³

Significand

M =	1. <u>1101101101₂</u>
frac=	<u>1101101101101</u> 000000000 ₂

Exponent

Ε	=	13	
Bias	=	127	
Ехр	=	140 =	100011002

Result:

0 10001100 1101101101000000000 s exp frac

Normalized Encoding Example 2

- Value: π, rounded to 24 bits of precision
 - sign: 0

Significand

S	=	11.0010010000111111011011 (including	hidden	bit)
М	=	1.10010010000111111011011 ₂ (x 2^1)		
frac	=	10010010000111111011011 ₂		

Exponent

Ε	=	1	
Bias	=	127	
Exp	=	128 =	10000000 ₂

Result:

0 1000000 1001000011111011011 s exp frac

Normalized Encoding Practice

Value: -π

same as before; only sign bit changes



Denormalized Values

- Also called denormal or subnormal numbers
- Values that are very close to zero
 - Fill the "underflow" gap around zero
 - Gradual underflow = numeric values are spaced evenly near 0.0
- Any number with magnitude smaller than the smallest normal number
 - When the exponent field is all zeros
 - E = 1-bias
 - Significand M = f without implied leading 1
 - h = 0 (hidden bit)

Representation of numeric value 0

 -0.0 and +0.0 are considered different in some ways and the same in others

Denormalized Values

- In a normal floating point value, there are no leading zeros in the significand, instead leading zeros are moved to the exponent.
- e.g., 0.0123 would be written as 1.23 * 10⁻²
- Denormalized numbers are numbers where this representation would result in an exponent that is too small (the exponent usually having a limited range). Such numbers are represented using leading zeros in the significand.

Denormalized Values

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

 $E = 1 - Bias$

Condition: exp = 000...0

Exponent value: E = 1 – Bias (instead of E = 0 – Bias)

- Significand coded with implied leading 0: M = 0.xxx...x₂
 - **xxx**...**x**: bits of **frac**

Cases

- **exp** = 000...0, **frac** = 000...0
 - Represents zero value
 - Note distinct values: +0 and -0 (why?)
- $exp = 000...0, frac \neq 000...0$
 - Numbers closest to 0.0
 - Equispaced

Denormalized Decoding Example

What is the decimal value of the following?

s exp frac

- sign: 0 (positive)
- we know it's a denormalized value because exp is all 0s

Exponent

E = 1 - 127 = -126 (same for all denormalized numbers)

Significand

frac =	0000001000011100000000 ₂
M =	$0.00000100011100000000_{2}$

Result:

- 0.0000001000011100000000 x $2^{-126} = 1.0000111 x 2^{-133}$
- 10000111 x 2^{-140} = 135 x 2^{-140} = 9.69 x 10^{-41}

Denormalized Encoding Example

Value: -25 15/32 x 2⁻¹³²

- sign: 1
- power of 2 indicates denormalized number (< -126)</p>

Exponent

Bias = 127 $Exp = 00000000_2$ (same for all denormalized numbers) E = 1 - 127 = -126

Significand

S	=	11001.01111 x 2⁻¹³² (move 6 places left to match -126)
М	=	0.01100101111000000000000000 ₂ (x 2 ⁻¹²⁶)
frac	=	01100101111000000000000000002

Result: 1 0000000 01100101110000000000 s exp frac

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 ≠ 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



Interesting Numbers

{single,double}

Description	ехр	frac	Numeric Value
Zero	0000	0000	0.0
Smallest Pos. Denorm.	0000	0001	2 ^{-{23,52}} x 2 ^{-{126,1022}}
Single ≈ 1.4 x 10 ⁻⁴⁵			
Double ≈ 4.9 x 10 ⁻³²⁴			
Largest Denormalized	0000	1111	$(1.0 - \varepsilon) \ge 2^{-\{126, 1022\}}$
Single ≈ 1.18 x 10 ⁻³⁸			
Double ≈ 2.2 x 10 ⁻³⁰⁸			
Smallest Pos. Normalized	0001	0000	1.0 x 2 ^{-{126,1022}}
Just larger than largest denorm	nalized		
One	0111	0000	1.0
Largest Normalized	1110	1111	(2.0 – ε) x 2 ^{127,1023}
Single ≈ 3.4 x 10 ³⁸			

Double $\approx 1.8 \times 10^{308}$

Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition

Examples and properties

- Rounding, addition, multiplication
- Floating point in C
- Summary

Tiny Floating Point Example

S	exp	frac
1	4-bits	3-bits

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

Normalize

S	exp	frac
1	4-bits	3-bits

Requirement

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

Value	Binary	Fraction	Exponent
128	1000000	1.0000000	7
13	00001101	1.1010000	3
17	00010001	1.0001000	4
19	00010011	1.0011000	4
138	10001010	1.0001010	7
63	00111111	1.1111100	5

Dyna	m	nic F	Range	e (Po	ositive	2	Only) [v = (-1) ^s M 2 ^E
	S	exp	frac	E	Value			n: E = Exp – Bias
	0	0000	000	-6	0			d: E = 1 – Bias
	0	0000	001	-6	1/8*1/64	=	1/512	closest to zero
Denormalized	0	0000	010	-6	2/8*1/64	=	2/512	
numbers	•••							
	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	
	0	0110	111	-1	15/8*1/2	=	15/16	closest to 1 below
Normalized	0	0111	000	0	8/8*1	=	1	
numbers	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			

Distribution of Values

6-bit IEEE-like format

- e = 3 exponent bits
- f = 2 fraction bits
- Bias is 2³⁻¹-1 = 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

S	exp	frac
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

Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

Floating Point Operations: Basic Idea

$$\mathbf{x} +_{f} \mathbf{y} = \text{Round}(\mathbf{x} + \mathbf{y})$$

x
$$\times_{f}$$
 y = 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 (– ∞)	\$1	\$1	\$1	\$2	-\$2
Round up (+ ∞)	\$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 underestimated

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...2

Examples

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

Value	Binary	Rounded	Action	Rounded Value
2 3/32	10.00 <mark>011</mark> 2	10.002	(<1/2—down)	2
2 3/16	10.00 <mark>110</mark> 2	10.012	(>1/2—up)	2 1/4
2 7/8	10.11 <mark>100</mark> 2	11.002	(1/2—up)	3
2 5/8	10.10 <mark>100</mark> 2	10.102	(1/2—down)	2 1/2



Round up conditions

- Round = 1, Sticky = 1 → > 0.5
- Guard = 1, Round = 1, Sticky = 0 → Round to even

Value	Fraction	GRS	Incr?	Rounded
128	1.0000000	000	N	1.000
13	1.1010000	100	N	1.101
17	1.0001000	010	N	1.000
19	1.0011000	110	Y	1.010
142	1.0001110	011	Y	1.001
63	1.111 <mark>1100</mark>	111	Y	10.000

Postnormalize

Issue

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

Value	Rounded	Ехр	Adjusted	Result
128	1.000	7		128
13	1.101	3		13
17	1.000	4		16
19	1.010	4		20
142	1.001	7		144
63	10.000	5	1.000/6	64

FP Multiplication

■ (-1)^{s1} M1 2^{E1} x (-1)^{s2} M2 2^{E2}

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

Sign s:	s1 ^ s2
Significand M:	M1 x M2
Exponent E:	E1 + E2

Fixing

- If $M \ge 2$, shift *M* right, increment *E*
- If E out of range, overflow
- Round *M* to fit **frac** precision

Implementation

Biggest chore is multiplying significands

FP Multiplication Example

What is the product of the following?

1 10011100	110000000000000000000000000000000000000	0000000
1 11110000	01100000000000000	0000000
s exp	frac	
Sign		1.110
1 ^1 = 0		<u>x 1.011</u>
		1 110
Exponent		11 100
<i>E1</i> = 156 - 127= 29 <i>E2</i>	2 = 240 - 127= 113	<u>+ 1110 000</u>
F = 29 + 113 + 1 = 143	+ 127 = 270 (1 0000 1110 - overflow)	10.011 010
		$= 1.001101 \times 2^{1}$

Significand

00

Result:

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

Get binary points lined up



Fixing

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

FP Addition Example

What is the sum of the following?



Sign

sign with larger exp = 0

Exponent

- *E1* = 125 127 = -2 *E2* = 129 127 = 2
- $E = 1000\ 0001\ (E2, \text{ or } 2 + 127 = 129)$

Significand

frac = **1101 0000 0000 0000 010**₂ (after round to even)

Result:

Bryant and O'Hallaron, Computer Systems: A Programmer's Perspective, Third Edition

x 2²

.01 0000 0000 0000 0000 0011 000

+ 111 .00 0000 0000 0000 0000 0000 000

 $= 1.1101\ 0000\ 0000\ 0000\ 0000\ 0011\ \times 2^{2}$

 $= 1.1101\ 0000\ 0000\ 0000\ 010$

111.010000000000000000011000

Floating Point

- Background: Fractional binary numbers
- IEEE floating point standard: Definition
- Example and properties
- Rounding, addition, multiplication
- Floating point in C
- Summary

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

• int \rightarrow float

- Cannot overflow; will round according to rounding mode
- int/float \rightarrow double
 - Exact conversion, as long as int has ≤ 53 bit word size
- float/double ightarrow int
 - Truncates fractional part; like rounding toward zero
 - Not defined when out of range or NaN: Generally sets to TMin
- double \rightarrow float
 - Can overflow (range smaller); may be rounded (precision smaller)

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

F	•	x == (int)(fl	oat) x
т	•	x == (int) (do	ouble) x
т	•	f == (float)(double) f
F	•	d == (double)	(float) d
т	•	f == -(-f);	
F	•	2/3 == 2/3.0	
т	•	d < 0.0	\Rightarrow ((d*2) < 0.0)
т	•	d > f =	$\Rightarrow -f > -d$
т	•	d * d >= 0.0	
F	•	(d+f)-d == f	

Summary

- IEEE Floating Point has clear mathematical properties
- Represents numbers of form M x 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