



# **Floating-point numbers and arithmetic**

Slides adapted by: Sparsh Mittal

# Quick Summary

$$(-1)^S \times (1.M) \times 2^E$$

		Range	Accuracy
FP32	<div> <div>1</div> <div>8</div> <div>23</div> <div>S</div> <div>E</div> <div>M</div> </div>	$10^{-38} - 10^{38}$	.000006%
FP16	<div> <div>1</div> <div>5</div> <div>10</div> <div>S</div> <div>E</div> <div>M</div> </div>	$6 \times 10^{-5} - 6 \times 10^4$	.05%
Int32	<div> <div>1</div> <div>31</div> <div>S</div> <div>M</div> </div>	$0 - 2 \times 10^9$	$\frac{1}{2}$
Int16	<div> <div>1</div> <div>15</div> <div>S</div> <div>M</div> </div>	$0 - 6 \times 10^4$	$\frac{1}{2}$
Int8	<div> <div>1</div> <div>7</div> <div>S</div> <div>M</div> </div>	$0 - 127$	$\frac{1}{2}$
Fixed point	<div> <div>S</div> <div>I</div> <div>F</div> </div> <div>↑ radix point</div>	-	-

Dally, High Performance Hardware for Machine Learning, NIPS'2015

# A Brief Guide to Floating Point Formats

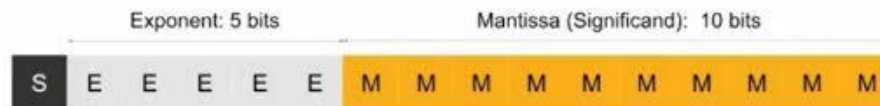
## fp32: Single-precision IEEE Floating Point Format

Range:  $\sim 1e^{-38}$  to  $\sim 3e^{38}$



## fp16: Half-precision IEEE Floating Point Format

Range:  $\sim 5.96e^{-8}$  to 65504



## bf16: Brain Floating Point Format

Range:  $\sim 1e^{-38}$  to  $\sim 3e^{38}$



# Floating Point

- Representation for non-integral numbers
  - Including very small and very large numbers
- Like scientific notation
  - $-2.34 \times 10^{56}$  ← normalized
  - $+0.002 \times 10^{-4}$  ← not normalized
  - $+987.02 \times 10^9$  ← not normalized
- In binary
  - $\pm 1.xxxxxxx_2 \times 2^{yyyy}$
- Types `float` and `double` in C

# Floating Point Standard

- Defined by IEEE Std 754-1985
- Developed in response to divergence of representations
  - Portability issues for scientific code
- Now almost universally adopted
- Two representations
  - Single precision (32-bit)
  - Double precision (64-bit)

# IEEE Floating-Point Format

single: 8 bits

double: 11 bits

single: 23 bits

double: 52 bits

S	Exponent	Fraction
---	----------	----------

$$x = (-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$$

- S: sign bit (0  $\Rightarrow$  non-negative, 1  $\Rightarrow$  negative)
- Normalize significand:  $1.0 \leq |\text{significand}| < 2.0$ 
  - Always has a leading pre-binary-point 1 bit, so no need to represent it explicitly (hidden bit)
  - Significand is Fraction with the “1.” restored
- Exponent: excess representation: actual exponent + Bias
  - Ensures exponent is unsigned
  - Single: Bias = 127; Double: Bias = 1023

# Single-Precision Range

- Exponents 00000000 and 11111111 reserved
- Smallest value
  - Exponent: 00000001  
 $\Rightarrow$  actual exponent =  $1 - 127 = -126$
  - Fraction: 000...00  $\Rightarrow$  significand = 1.0
  - $\pm 1.0 \times 2^{-126} \approx \pm 1.2 \times 10^{-38}$
- Largest value
  - exponent: 11111110  
 $\Rightarrow$  actual exponent =  $254 - 127 = +127$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$
  - $\pm 2.0 \times 2^{+127} \approx \pm 3.4 \times 10^{+38}$

# Double-Precision Range

- Exponents 0000...00 and 1111...11 reserved
- Smallest value
  - Exponent: 000000000001  
 $\Rightarrow$  actual exponent =  $1 - 1023 = -1022$
  - Fraction: 000...00  $\Rightarrow$  significand = 1.0
  - $\pm 1.0 \times 2^{-1022} \approx \pm 2.2 \times 10^{-308}$
- Largest value
  - Exponent: 111111111110  
 $\Rightarrow$  actual exponent =  $2046 - 1023 = +1023$
  - Fraction: 111...11  $\Rightarrow$  significand  $\approx 2.0$
  - $\pm 2.0 \times 2^{+1023} \approx \pm 1.8 \times 10^{+308}$



# Why it is called single/double precision

- The *precision* indicates the number of decimal digits that are **correct**, i.e. without any kind of representation error or approximation. In other words, it indicates how many decimal digits one can **safely** use.
- The number of decimal digits which can be safely used:
- **single precision**:  $\log_{10}(2^{24})$ , which is about 7~8 decimal digits
- **double precision**:  $\log_{10}(2^{53})$ , which is about 15~16 decimal digits

- A floating-point variable can represent a wider range of numbers than a fixed-point variable of the same bit width at the cost of precision.
- A signed 32-bit integer variable has a maximum value of  $2^{31} - 1 = 2,147,483,647$ , whereas an IEEE 754 32-bit base-2 floating-point variable has a maximum value of  $(2 - 2^{-23}) \times 2^{127} \approx 3.4028235 \times 10^{38}$ .

# Floating-Point Example

- Represent  $-0.75$ 
  - $-0.75 = (-1)^1 \times 1.1_2 \times 2^{-1}$
  - $S = 1$
  - Fraction =  $1000\dots00_2$
  - Exponent =  $-1 + \text{Bias}$ 
    - Single:  $-1 + 127 = 126 = 01111110_2$
    - Double:  $-1 + 1023 = 1022 = 01111111110_2$
- Single:  $10111111101000\dots00$
- Double:  $10111111111101000\dots00$

# Floating-Point Example

- What number is represented by the single-precision float

11000000101000...00

- $S = 1$
  - Fraction =  $01000...00_2$
  - Exponent =  $10000001_2 = 129$
- $x = (-1)^1 \times (1 + 01_2) \times 2^{(129 - 127)}$   
 $= (-1) \times 1.25 \times 2^2$   
 $= -5.0$

# Denormal Numbers

- Exponent = 000...0  $\Rightarrow$  hidden bit is 0

$$x = (-1)^s \times (0 + \text{Fraction}) \times 2^{-126}$$

- Smaller than normal numbers
  - allow for gradual underflow, with diminishing precision

# Denormal Numbers

- \* Smallest +ve normal number :  $2^{-126}$
- \* Largest denormal number :
  - \*  $0.11...11 * 2^{-126} = (1 - 2^{-23}) * 2^{-126}$   
 $* = 2^{-126} - 2^{-149}$

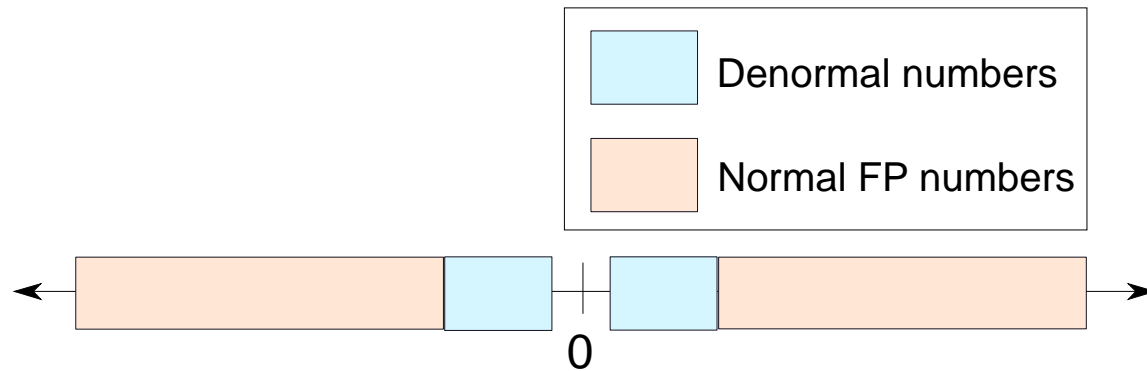
# Example

Find the ranges of denormal numbers.

## **Answer**

- For positive denormal numbers, the range is  $[2^{-149}, 2^{-126} - 2^{-149}]$
- For negative denormal numbers, the range is  $[-2^{-149}, -2^{-126} + 2^{-149}]$

# Denormal Numbers on Number Line





# Infinites and NaNs

- Exponent = 111...1, Fraction = 000...0
  - $\pm$ Infinity
  - Can be used in subsequent calculations, avoiding need for overflow check
- Exponent = 111...1, Fraction  $\neq$  000...0
  - Not-a-Number (NaN)
  - Indicates illegal or undefined result
    - e.g., 0.0 / 0.0
  - Can be used in subsequent calculations

# Special FP Numbers

$E$	$M$	Value
255	0	$\infty$ if $S = 0$
255	0	$-\infty$ if $S = 1$
255	$\neq 0$	NAN(Not a number)
0	0	0
0	$\neq 0$	Denormal number

- \*  $\text{NAN} + x = \text{NAN}$        $1/0 = \infty$
- \*  $0/0 = \text{NAN}$        $-1/0 = -\infty$
- \*  $\sin^{-1}(5) = \text{NAN}$

# **Some possible pitfalls in use of FP arithmetic**

# Associativity

- Parallel programs may interleave operations in unexpected orders
  - Assumptions of associativity may fail

		$(x+y)+z$	$x+(y+z)$
x	-1.50E+38	0.00E+00	-1.50E+38
y	1.50E+38		1.50E+38
z	1.0	1.0	
		1.00E+00	0.00E+00

- Need to validate parallel programs under varying degrees of parallelism

# Right Shift and Division

- Left shift by  $i$  places multiplies an integer by  $2^i$
- Right shift divides by  $2^i$ ?
  - Only for unsigned integers
- For signed integers
  - Arithmetic right shift: replicate the sign bit
  - e.g.,  $-5 / 4$ 
    - $11111011_2 \gg 2 = 11111110_2 = -2$
    - Rounds toward  $-\infty$
  - c.f.  $11111011_2 \ggg 2 = 00111110_2 = +62$



# **FP addition and multiplication**

# FP Addition (Base 10)

- Consider a 4-digit decimal example
  - $9.999 \times 10^1 + 1.610 \times 10^{-1}$
- 1. Align decimal points
  - Shift number with smaller exponent
  - $9.999 \times 10^1 + 0.016 \times 10^1$
- 2. Add significands
  - $9.999 \times 10^1 + 0.016 \times 10^1 = 10.015 \times 10^1$
- 3. Normalize result & check for over/underflow
  - $1.0015 \times 10^2$
- 4. Round and renormalize if necessary
  - $1.002 \times 10^2$

# FP Addition (Base 2)

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2}$
  - Above is same as  $(0.5 + -0.4375)$
- 1. Align binary points
  - Shift number with smaller exponent
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$
- 2. Add significands
  - $1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$
- 3. Normalize result & check for over/underflow
  - $1.000_2 \times 2^{-4}$ , with no over/underflow
- 4. Round and renormalize if necessary
  - $1.000_2 \times 2^{-4}$  (no change) = 0.0625



# Floating-Point Multiplication

- Consider a 4-digit decimal example
  - $1.110 \times 10^{10} \times 9.200 \times 10^{-5}$
- 1. Add exponents
  - For biased exponents, subtract bias from sum
  - New exponent =  $10 + -5 = 5$
- 2. Multiply significands
  - $1.110 \times 9.200 = 10.212 \Rightarrow 10.212 \times 10^5$
- 3. Normalize result & check for over/underflow
  - $1.0212 \times 10^6$
- 4. Round and renormalize if necessary
  - $1.021 \times 10^6$
- 5. Determine sign of result from signs of operands
  - $+1.021 \times 10^6$

# Floating-Point Multiplication

- Now consider a 4-digit binary example
  - $1.000_2 \times 2^{-1} \times -1.110_2 \times 2^{-2}$  ( $0.5 \times -0.4375$ )
- 1. Add exponents
  - Unbiased:  $-1 + -2 = -3$
  - Biased:  $(-1 + 127) + (-2 + 127) = -3 + 254 - 127 = -3 + 127$
- 2. Multiply significands
  - $1.000_2 \times 1.110_2 = 1.1102 \Rightarrow 1.110_2 \times 2^{-3}$
- 3. Normalize result & check for over/underflow
  - $1.110_2 \times 2^{-3}$  (no change) with no over/underflow
- 4. Round and renormalize if necessary
  - $1.110_2 \times 2^{-3}$  (no change)
- 5. Determine sign:  $+ve \times -ve \Rightarrow -ve$ 
  - $-1.110_2 \times 2^{-3} = -0.21875$

# Further study

- <https://blog.demofox.org/2017/11/21/floating-point-precision/>
- <https://stackoverflow.com/questions/4220417/print-binary-representation-of-a-float-number-in-c>