



# COMPSCI 210

## Part III

### Floating Point Numbers

## Agenda & Reading

### Agenda:

- Introduction
- Normalization
- IEEE 754 Floating Point Representation
  - Single Precision (32-bit)
  - Double Precision (64-bit)
- Examples
- Special Numbers
- Conversion

### Recommended Reading:

- IEEE Floating Point Numbers
  - [http://en.wikipedia.org/wiki/IEEE\\_floating-point\\_standard](http://en.wikipedia.org/wiki/IEEE_floating-point_standard)

### Animation

- Float.htm

### Exercise

- Worksheet 02

III-02-2

## Introduction

Given the following numbers:

- $37.25_{10} = 100101.01_2$
- $0.3125_{10} = 0.0101_2$
- $+300000000 \text{ ms}^{-1}$
- $+299,792,458 \text{ ms}^{-1}$

How can we store the above number in 4 bytes?

People usually represent very large and small numbers in “scientific notation”, as a fixed point number times a power of 10.

### Scientific Notation:

- $3 * 10^8$
- $2.9979 * 10^8$
- $3.125 * 10^{-1}$
- $3.725 * 10^1$

### Floating point numbers are represented in the computer in a similar manner, but using base 2 rather than base 10.

- $1.0010101 * 2^5$
- $1.11101 * 2^2$
- $1.01 * 2^{-2}$

III-02-3

## Normalization

Given the following number

- $12.34 = 0.1234 * 10^2 = 1.234 * 10^1 = 123.4 * 10^{-1}$

Scientific numbers are always written with one digit before the point, giving a “normalized” representation.

- After Normalization =  $1.234 * 10^1$

Normalization in Base 2

- One digit to the left of the binary point. It must be 1.

### Examples:

➢  $100101.01_2 = 1.0010101 * 2^5$  Radix point move to left by 5 places

➢  $111.101_2 = 1.11101 * 2^2$  Radix point move to left by 2 places

➢  $0.0101_2 = 1.01 * 2^{-2}$  Radix point move to right by 2 places

After normalization, the numbers now have a standard format

III-02-4



## Special Numbers

### Special Exponents (Single)

- 00000000
  - Case 1: Represent number in Denormalized format!
    - Exponent is all zeros, the floating-point number is Denormalized
    - =  $(0.0 + \text{significand}) * 2^{-126}$
  - Case 2: Represent ZERO with all zeros in exponent and mantissa bits
    - +0, -0
- 11111111
  - NaN
    - An exponent of all ones with any other mantissa is interpreted to mean "not a number". Example: 0 11111111 000000000000000000000001
  - Infinity:
    - An exponent of all ones with a mantissa whose bits are all zero indicates an infinity. The sign of the infinity is indicated by the sign bit.
    - Example: 0 11111111 000000000000000000000000

III-02-9

## Table

Range Name		
-NaN	1 11...11 00...01	FF800001
-Infinity	1 11...11 00...00	FF800000
Negative Normalized	1 11..10 11..11	FF7FFFFF
	...	
Negative Normalized	1 00...01 00..00	80800000
Negative Denormalized	1 00..00 11...11	807FFFFF
	...	
Negative Denormalized	1 00..00 00...01	80000001
-0	1 00..00 00..00	80000000
+0	0 00..00 00..00	00000000
Positive Denormalized	0 00..00 00...01	00000001
	...	
Positive Denormalized	0 00..00 11...11	007FFFFF
Positive Normalized	0 00..01 00...00	00800000
	...	
Positive Normalized	0 11..10 11...11	7F7FFFFF
+Infinity	0 11..11 00..00	7F800000
+NaN	0 11..11 00..01	7F800001

II-02-10

## Conversion:

Example:  
IEEE754.html (Enter a float value, select the type (Single/double))

### Example 1: Decimal -> IEEE Floating Point

- $-1.25_{10}$ ,
  - Binary = -1.01, Normalization  $\Rightarrow -1.01 = -1.01 * 2^0$
  - Answer:
    - Exponent = 127 = 01111111
    - Sign = 1
    - Mantissa = 010...0
  - Answer = 1 01111111 010...0 = BFA00000

Shockwave Movie: [float.htm](#)

### Example 2: IEEE Floating Point -> Decimal

- 40900000
  - = 0100 0000 1001 0000 ... 0000
  - Answer:
    - Sign = 0
    - Exponent = 100 0000 1 = 129 = 129-127
    - Mantissa = 001 0000 ... 0000 = 0.125
  - Answer =  $(-1)^0 * (1.0 + 0.125) * 2^{(2)} = 4.5$

III-02-11

## Java/C Examples

### Print Hexadecimal Strings (Java)

```
Integer.toHexString(Float.floatToIntBits (f))
Long.toHexString(Double.doubleToLongBits (d))
```

### Print decimal value (Java)

```
Float.intBitsToFloat(0xbfa00000)
Double.longBitsToDouble(0x4012000000000000L)
```

### Print Hexadecimal Strings (C)

```
float num = -1.25;
int temp;
temp = *((int*) &num); //convert the float to the Hex bit patterns
printf("\nThe number is %f = %x", num, temp);
```

III-02-12

## Range & Problems

### Range:



- Magnitude of numbers that can be represented is in the range:
- Single Precision (Normalized)
  - $2^{-126} * (1.0)$  to  $2^{127} * (2 - 2^{-23})$  Why? (see next page)
    - $2^{-126} * (1.0)$ : Exponent: 00000001 Mantissa: 000...000
    - $2^{127} * (2 - 2^{-23})$ : Exponent: 11111110, Mantissa: 1...111 =  $3.40 \times 10^{38}$
    - $1.18 \times 10^{-38}$  to  $3.40 \times 10^{38}$  (approximately)
- Double Precision
  - $2^{-1022} * (1.0)$  to  $2^{1023} * (2 - 2^{-52})$  ( $2.23 * 10^{-308}$  to  $1.8 * 10^{308}$ )

### Problems

- Overflow occurs when the exponent is larger than the allocated space
- Underflow occurs when a negative exponent is too large in absolute value to fit within the bits allocated to store it
- Truncation occurs when there are not enough digits in the mantissa to represent the number:
  - Examples: 1/3, 1/10 etc

$1/10 = 3DCCCCC = (1.1001100110011...) \times 2^{-4} = 0.09999...$   
 $1/3 = 3EAAAAA = (1.010101010101...) \times 2^{-2} = 0.3333...$

III-02-13

## Extra

### 0 00..01 00...00

- $= (1.0 + m) * 2^e$
- $= (1.0 + 0) * 2^{1-127}$
- $= 2^{-126}$

### 0 11..10 11...11

- $= (1.0 + m) * 2^e$
- $= (1.0 + 2^{-1} + 2^{-2} + 2^{-3} + \dots + 2^{-23}) * 2^{254-127}$
- $= (1.0 + 2^{-23} (2^{22} + \dots + 2^2 + 2 + 1)) * 2^{127}$
- $= (1.0 + 2^{-23} (1 + 2 + 2^2 + \dots + 2^{22})) * 2^{127}$
- $= (1.0 + 2^{-23} (2^{23} - 1) / (2 - 1)) * 2^{127}$
- $= (1.0 + 2^{-23} (2^{23} - 1)) * 2^{127}$
- $= (1.0 + 2^0 - 2^{-23}) * 2^{127}$
- $= (2 - 2^{-23}) * 2^{127}$

Geometric series:  
 $1 + x + x^2 + \dots + x^n$   
 $= \frac{x^{n+1} - 1}{x - 1}$

III-02-14