# RUTGERS

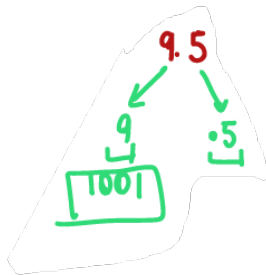## CS211 Computer Architecture
### Fall 2020

Recitation 6

# Floating Point Representation

- Used to represent more precise values (decimals)
- In order to do this, we must be able to know how to convert fractions into binary

# Fractions in Binary

Ex: 9.5

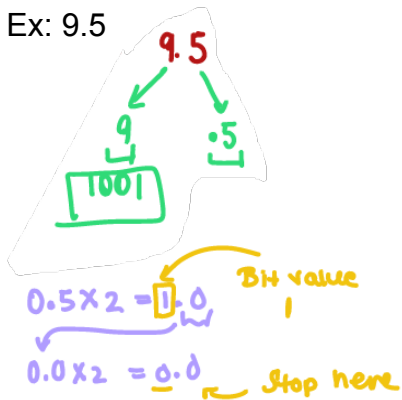- Separate the values before and after the decimal and convert

9.5

9 → .5

$\boxed{1001}$

if numbers
before decimal place is $2^0$, what about
<u>after</u> ?

$2^{-1}, 2^{-2}, \dots \quad 2^{-n}$

# Fractions in Binary

Ex: 9.5



to convert whole numbers
into binary, we divided.
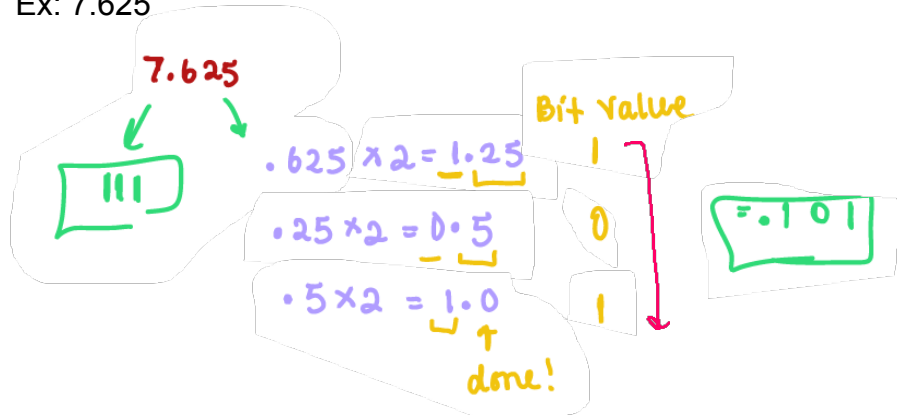For decimals, we multiply.

Multiply the decimals by 2, the
corresponding bit value is the number to
the left of the decimal place

Note: we do NOT reverse the order of
bits as we did for whole numbers. So
we are left with **0b.10** and you can drop
the 0s after the 1 to get **0b0.1**

$0.5 \times 2 = \boxed{1}.0$   Bit value

$0.0 \times 2 = 0.0$   stop here

$0.5 = \boxed{.1}$

# Fractions in Binary

Ex: 7.625

7.625

↙    ↘

111

.625 × 2 = 1.25

.25 × 2 = 0.5

.5 × 2 = 1.0
↑
done!

Bit value

1

0

1

= .101

# Fractions in Binary

Ex: 7.625

$$\begin{array}{ccc} -1 & -2 & -3 \\ =.1 & 0 & 1 \end{array}$$

<u>VERIFY</u>

$2^{-1} + 0 \cdot 2^{-2} + 2^{-3}$

$= \frac{1}{2} + \frac{1}{8} = .5 + .125 = .625$

7.625 = **0b111.101**

# IEEE Floating Point

- Let's focus on **single precision** – 32 bits
  - Normalized

- Anatomy of Floating Point

  - sign bit – 1
  - exponent bias – 8
  - mantissa – 23
  } 32

# Steps for General Floating Point Rep

General Steps:

① notice sign bit, then ignore it

② convert to binary (magnitude)

③ move decimal point before first 1 bit (keep track of how many spots you move)

④ identify mantissa as everything after decimal

⑤ Bias = $2^{n-1} - 1$    (n will be given-usually 8)

⑥ Put it all together

# RUTGERS

## Example

$4/2 = 2$  0
$2/2 = 1$  0
$1/2 = 0$  1

$100.11$

$.75/2$
$1.5 > 1$

$1.001 \times b^2$

$0\ |\underbrace{0000001}\ 0011 \longrightarrow$

$129$

$128/2 = 64$ 1
$64/2 = 32$ 0
$32/2 = 16$ 0
$16/2 = 8$ 0
$8/2 = 4$ 0
$4/2 = 2$ 0
$2/2 = 1$ 0

Ex: -4.75 in IEEE single-precision

① negative = 1

② $4.75 = 100.11$

$.75 \times 2 = 1.50$   1

$.50 \times 2 = 1.0$   1

Bit

# Example

Ex: -4.75 in IEEE single-precision

③ $100.11 \rightarrow 1.0011 \times 2^{②}$ ← moved 2 places

④ mantissa = 0011

# Example

Ex: -4.75 in IEEE ~~single precision~~ Single precision

Bias is 8 bit as we use 8 bits to represent the exponent

⑤ exponent = bias + decimal places moved

Bias, n=8

$2^8 - 1 = 127$

$127 + (2) = \underline{129}$

convert to Binary in 8 bits (n) $= 10000001$
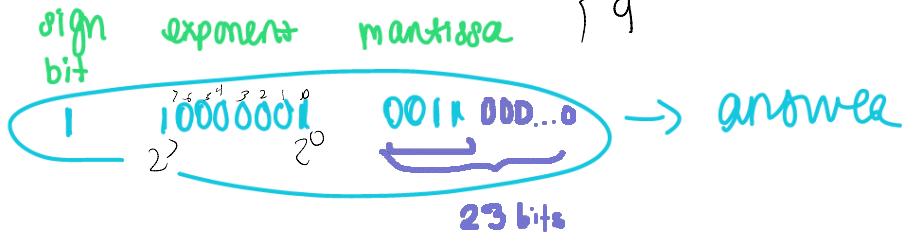
# Example

Ex: -4.75 in IEEE single-precision



b

sign bit | exponent | mantissa

$23$

$\frac{4}{19}$

| 7 6 5 4 3 2 1 0 | |
|---|---|
| 1 | 10000001 | 0011 000...0 | → answer |

$2^7$    $2^0$

23 bits

# Example

Ex: -4.75 in IEEE single-precision

how do you check? work backwards ☺

- exponent -127 = decimal places moved = e
- 1.[mantissa] x $2^e$    0.4
- sign = 1 if negative
         0 if positive

# Recap

- what does IEEE floating point mean?
  754

  different types of representation
  - single pt — 32 bit
  - double pt — 64 bit

  ↳ 32 bits in answer

  ↳ 1 sign bit

  ↳ 8 exponent

  ↳ 32-1-8 = 23 for mantissa (just add 0s)

## Quick Note - Normalized vs Denormalized

- The range of values you can represent in floating point can be defined as **normalized** or **denormalized**
- **Normalized**
  - Our exponent bias = <u>exponent – bias</u>
  - Value represented leading with <u>1 (1.[mantissa])</u>
- **Denormalized**
  - Our exponent bias = <u>1 – bias (</u>exponent = 0)
  - Value represented leading with <u>0 (0.[mantissa])</u>
  - Usually represents numbers very close to 0