

# CSC 411

Computer Organization (Spring 2024)  
Lecture 7: Floating Point

Prof. Marco Alvarez, University of Rhode Island

## Fractional binary numbers

- Bits to the right of binary point
  - fractional powers of 2
  - don't worry about negatives for now

<sup>2<sup>i</sup></sup>	<sup>2<sup>i-1</sup></sup>		<sup>2</sup>	<sup>1</sup>	<sup>1/2</sup>	<sup>1/4</sup>		<sup>2<sup>-j+1</sup></sup>	<sup>2<sup>-j</sup></sup>
b <sub>i</sub>	b <sub>i-1</sub>	...	b <sub>1</sub>	b <sub>0</sub>	b <sub>-1</sub>	b <sub>-2</sub>	...	b <sub>-j+1</sub>	b <sub>-j</sub>

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

$$11.010 =$$

$$1101.101 =$$

## Practice

- Convert fractional binary numbers to decimal

$$1.10 =$$

$$11.001 =$$

$$100.01 =$$

$$11.111 =$$

## Observations

- Not all decimal fractions have exact binary equivalents

• can only represent numbers of the form  $\frac{x}{2^k}$

- e.g., 1/5 and 1/10

- Limited precision due to finite number of bits

- can't easily represent very small or very large values

- 0.111111... represents a number just below 1.0

$$\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots + \frac{1}{2^i} + \dots = 1.0 - \epsilon$$

## IEEE Floating Point

### ▸ IEEE standard 754

- defines a common format for representing real numbers in computers
  - developed in response to divergence of representations
- supported by major CPUs (almost universally adopted)
- provides different precision levels (single, double, extended) for various needs.
- standardizes sign, exponent, and fraction components



## Floating point representation

### ▸ Numerical form

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

- **sign bit s**: 0 for positive, 1 for negative
- **exponent E**: magnitude of the number (power of 2)
  - encoded in exp
- **significant M**: captures the fractional part, scaled by the exponent, normally in range [1.0, 2.0)
  - encoded in frac



## Precision options

### ▸ Single-precision (32 bits)

- good balance of performance and range (7 decimal digits)



### ▸ Double-precision (64 bits)

- higher precision (15-17 decimal digits) for demanding calculations



### ▸ Others

- half precision, quad precision

## Normalized and denormalized numbers

### ▸ Normalized

- maximizes precision
- **exp** != 000...000 and **exp** != 111...111

### ▸ Denormalized

- used for very small numbers, reducing precision
- **exp** == 000...000

### ▸ Special

- **exp** == 111...111

