#### CSC 252/452: Computer Organization Fall 2024: Lecture 3

Instructor: Yanan Guo Department of Computer Science University of Rochester

## Announcement

- Programming Assignment 1 is out
  - Details:

https://www.cs.rochester.edu/courses/252/fall2024/labs/ assignment1.html

- Due on Sep 16th, 11:59 PM
- You have 3 slip days

## Announcement

- Programming Assignment 1 is in C language.
- Seek help from TAs.
  - TAs are best positioned to answer your questions about programming assignments!!!
- Programming assignments do NOT repeat the lecture materials. They ask you to synthesize what you have learned from the lectures and work out something new.
- Pay attention to Blackboard announcements
  - There are changes about the office hour locations/time...
  - I have to move my office hour tomorrow to early next week.

#### Last Lecture

- Why Binary (bits)?
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary

Two's Complement

• Two's Complement



| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

• Two's Complement

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

Two's Complement

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

 $b_2b_1b_0$ 

22 21 20

Weights in Unsigned

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

Weights in Unsigned

Weights in

Signed

2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

-2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

 $b_2b_1b_0$ 

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

 $b_2b_1b_0$ 

↗↑₹

Weights in Unsigned 2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

Weights in **-2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>** Signed

 $101_2 = 1^*2^0 + 0^*2^1 + (-1^*2^2) = -3_{10}$ 

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

Weights in Unsigned

Weights in Signed -2<sup>2</sup>

**-2**<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

 $b_2b_1b_0$ 

2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

$$101_2 = 1^* 2^0 + 0^* 2^1 + (-1^* 2^2) = -3_{10}$$

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

Weights in Unsigned

Weights in Signed -2<sup>2</sup>

**-2**<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

 $b_2b_1b_0$ 

2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

$$101_2 = 1^* 2^0 + 0^* 2^1 + (-1^* 2^2) = -3_{10}$$

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

Two's Complement

-4 -3 -2 -1 0 1 2 3

Weights in Unsigned

Weights in Signed -2<sup>2</sup>

**-2**<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

 $b_2b_1b_0$ 

2<sup>2</sup> 2<sup>1</sup> 2<sup>0</sup>

$$101_2 = 1^* 2^0 + 0^* 2^1 + (-1^* 2^2) = -3_{10}$$

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

- Only 1 zero
- There is (still) a bit that represents sign!
- Unsigned arithmetic still works

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

| • | Only | 1 zero |
|---|------|--------|
|---|------|--------|

- There is (still) a bit that represents sign!
- Unsigned arithmetic still works

010 +) 101 111

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

| • | Only | 1 | zero |
|---|------|---|------|
|---|------|---|------|

- There is (still) a bit that represents sign!
- Unsigned arithmetic still works

|    | 010 | 2     |
|----|-----|-------|
| +) | 101 | +) -3 |
|    | 111 | -1    |

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

| • Only 1 zero                     | )                 |                  | Signed | Binary |
|-----------------------------------|-------------------|------------------|--------|--------|
| • There is (at                    |                   | onvocanto signal | 0      | 000    |
| There is (st                      | iii) a bit that r | epresents sign:  | 1      | 001    |
| • Unsigned arithmetic still works |                   |                  | 2      | 010    |
| -                                 |                   |                  | 3      | 011    |
|                                   |                   |                  | -4     | 100    |
| 01                                | 0                 | 2                | -3     | 101    |
| +) 10                             | 1                 | +) -3            | -2     | 110    |
| 11                                | 1                 | -1               | -1     | 111    |

• 3 + 1 becomes -4 (called overflow. More on it later.)

 Suppose you want to define a variable that represents a person's age. What assumptions can you make about this variable's numerical value?

- Suppose you want to define a variable that represents a person's age. What assumptions can you make about this variable's numerical value?
  - Integer
  - Non-negative
  - Between 0 and 255 (8 bits)

- Suppose you want to define a variable that represents a person's age. What assumptions can you make about this variable's numerical value?
  - Integer
  - Non-negative
  - Between 0 and 255 (8 bits)
- Define a data type that captures all these attributes: unsigned char in C
  - Internally, an **unsigned char** variable is represented as a 8-bit, non-negative, binary number

• What if you want to define a variable that could take negative values?

- What if you want to define a variable that could take negative values?
  - That's what signed data types (e.g., **int**, **short**, etc.) are for

- What if you want to define a variable that could take negative values?
  - That's what signed data types (e.g., **int**, **short**, etc.) are for
- How are int values internally represented?
  - Theoretically could be either signed-magnitude or two's complement
  - The C language designers chose two's complement

|      |      | W       |                |                            |  |
|------|------|---------|----------------|----------------------------|--|
|      | 8    | 16      | 32             | 64                         |  |
| UMax | 255  | 65,535  | 4,294,967,295  | 18,446,744,073,709,551,615 |  |
| TMax | 127  | 32,767  | 2,147,483,647  | 9,223,372,036,854,775,807  |  |
| TMin | -128 | -32,768 | -2,147,483,648 | -9,223,372,036,854,775,808 |  |

| C Data Type  |      | 32-bit | 64-bit |
|--------------|------|--------|--------|
| (unsigned) c | har  | 1      | 1      |
| (unsigned) s | hort | 2      | 2      |
| (unsigned) i | nt   | 4      | 4      |
| (unsigned) 1 | ong  | 4      | 8      |

|      |      | W       |                |                            |  |
|------|------|---------|----------------|----------------------------|--|
|      | 8    | 16      | 32             | 64                         |  |
| UMax | 255  | 65,535  | 4,294,967,295  | 18,446,744,073,709,551,615 |  |
| TMax | 127  | 32,767  | 2,147,483,647  | 9,223,372,036,854,775,807  |  |
| TMin | -128 | -32,768 | -2,147,483,648 | -9,223,372,036,854,775,808 |  |

| C Data Type      | 32-bit | 64-bit |
|------------------|--------|--------|
| (unsigned) char  | 1      | 1      |
| (unsigned) short | 2      | 2      |
| (unsigned) int   | 4      | 4      |
| (unsigned) long  | 4      | 8      |

- C Language
  - •#include <limits.h>
  - •Declares constants, e.g.,

• ULONG MAX

- LONG MAX
- LONG\_MIN

Values platform specific

# Mapping Between Signed & Unsigned

 Mappings between unsigned and two's complement numbers: Keep bit representations and reinterpret

| Signed | Unsigned | Binary |
|--------|----------|--------|
| 0      | 0        | 000    |
| 1      | 1        | 001    |
| 2      | 2        | 010    |
| 3      | 3        | 011    |
| -4     | 4        | 100    |
| -3     | 5        | 101    |
| -2     | 6        | 110    |
| -1     | 7        | 111    |

## Mapping Signed $\leftrightarrow$ Unsigned

| Bits | Signed |                       |
|------|--------|-----------------------|
| 0000 | 0      |                       |
| 0001 | 1      |                       |
| 0010 | 2      |                       |
| 0011 | 3      |                       |
| 0100 | 4      |                       |
| 0101 | 5      | —→T2U—→               |
| 0110 | 6      |                       |
| 0111 | 7      |                       |
| 1000 | -8     | ← <u>U2I</u> ←        |
| 1001 | -7     |                       |
| 1010 | -6     | +/- 16                |
| 1011 | -5     | $\longleftrightarrow$ |
| 1100 | -4     |                       |
| 1101 | -3     |                       |
| 1110 | -2     |                       |
| 1111 | -1     |                       |



#### **Today: Representing Information in Binary**

- Why Binary (bits)?
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary

#### **The Problem**

| short | int x = | = | 15213;   |
|-------|---------|---|----------|
| int   | ix =    | = | (int) x; |
| short | int y = | = | -15213;  |
| int   | iy =    | = | (int) y; |

| C Data Type | 64-bit |
|-------------|--------|
| char        | 1      |
| short       | 2      |
| int         | 4      |
| long        | 8      |

#### **The Problem**

short int x = 15213; int ix = (int) x; short int y = -15213; int iy = (int) y;

| C Data Type | 64-bit |
|-------------|--------|
| char        | 1      |
| short       | 2      |
| int         | 4      |
| long        | 8      |

- Converting from smaller to larger integer data type
- Should we preserve the value?
- Can we preserve the value?
- How?

#### **The Problem**

short int x = 15213; int ix = (int) x; short int y = -15213; int iy = (int) y;

| C Data Type | 64-bit |
|-------------|--------|
| char        | 1      |
| short       | 2      |
| int         | 4      |
| long        | 8      |

- Converting from smaller to larger integer data type
- Should we preserve the value?
- Can we preserve the value?
- How?

|    | Decimal | Hex         | Binary                             |
|----|---------|-------------|------------------------------------|
| x  | 15213   | 3B 6D       | 00111011 01101101                  |
| ix | 15213   | 00 00 3B 6D | 0000000 0000000 00111011 01101101  |
| У  | -15213  | C4 93       | 11000100 10010011                  |
| iy | -15213  | FF FF C4 93 | 11111111 1111111 11000100 10010011 |

# **Signed Extension**

- Task:
  - Given *w*-bit signed integer *x*
  - Convert it to (w+k)-bit integer with same value

# **Signed Extension**

- Task:
  - Given *w*-bit signed integer *x*
  - Convert it to (*w*+*k*)-bit integer with same value
- Rule:
  - Make *k* copies of sign bit:

• 
$$X' = x_{w-1}, ..., x_{w-1}, x_{w-1}, x_{w-2}, ..., x_0$$

k copies of MSB

# **Signed Extension**

- Task:
  - Given *w*-bit signed integer *x*
  - Convert it to (w+k)-bit integer with same value
- Rule:
  - Make *k* copies of sign bit:

• 
$$X' = x_{w-1}, ..., x_{w-1}, x_{w-1}, x_{w-2}, ..., x_0$$



#### **Another Problem**

unsigned short x = 47981; unsigned int ux = x;

|    | Decimal | Hex         | Binary                              |
|----|---------|-------------|-------------------------------------|
| x  | 47981   | BB 6D       | 10111011 01101101                   |
| ux | 47981   | 00 00 BB 6D | 00000000 00000000 10111011 01101101 |
# **Unsigned (Zero) Extension**

- Task:
  - Given *w*-bit unsigned integer *x*
  - Convert it to (w+k)-bit integer with same value
- Rule:
  - Simply pad zeros:

• 
$$X' = 0, ..., 0, x_{w-1}, x_{w-2}, ..., x_0$$



#### Yet Another Problem

int x = 53191; short sx = (short) x;

|    | Decimal | Hex         | Binary                            |
|----|---------|-------------|-----------------------------------|
| x  | 53191   | 00 00 CF C7 | 0000000 0000000 11001111 11000111 |
| sx | -12345  | CF C7       | 11001111 11000111                 |

#### **Yet Another Problem**

| int   | х  | = | 53191;  | 1; |  |
|-------|----|---|---------|----|--|
| short | sx | = | (short) | x  |  |

|    | Decimal | Hex         | Binary                            |
|----|---------|-------------|-----------------------------------|
| x  | 53191   | 00 00 CF C7 | 0000000 0000000 11001111 11000111 |
| sx | -12345  | CF C7       | 11001111 11000111                 |

- Truncating (e.g., int to short OR unsigned int to unsigned short)
  - C's implementation: leading bits are truncated, results reinterpreted
  - So can't always preserve the numerical value

#### **Today: Representing Information in Binary**

- Why Binary (bits)?
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary
- Representations in memory, pointers, strings

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

• Similar to Decimal Addition

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

- Similar to Decimal Addition
- Suppose we have a new data type that is 3-bit wide (c.f., short has 16 bits)

|        | 010    | 2    |  |
|--------|--------|------|--|
| Normal | +) 101 | +) 5 |  |
| Case   | 111    | 7    |  |

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

- Similar to Decimal Addition
- Suppose we have a new data type that is 3-bit wide (c.f., short has 16 bits)
- Might overflow: result can't be represented within the size of the data type

| Normal   | 010<br>+) 101 | 2<br>+) 5 |
|----------|---------------|-----------|
| Case     | 111           | 7         |
| Overflow | 110<br>+) 101 | 6<br>+) 5 |
| Case     | 1011          | 11        |

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

- Similar to Decimal Addition
- Suppose we have a new data type that is 3-bit wide (c.f., short has 16 bits)
- Might overflow: result can't be represented within the size of the data type

| Normal   | 010<br>+) 101 | 2<br>+) 5 | 6<br>7   |
|----------|---------------|-----------|----------|
| Case     | 111           | 7         |          |
| Overflow | 110<br>+) 101 | 6<br>+) 5 |          |
| Case     | 1011          | 11        | True Sum |

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

- Similar to Decimal Addition
- Suppose we have a new data type that is 3-bit wide (c.f., short has 16 bits)
- Might overflow: result can't be represented within the size of the data type

| Normal<br>Case | 010           | 2<br>+) 5        | 6<br>7                                      | 110<br>111            |
|----------------|---------------|------------------|---|-----------------------|
|                | 111           | <u>+) 5</u><br>7 |   |                       |
| Overflow       | 110<br>+) 101 | 6<br>+) 5        |   |                       |
| Case           | 1011<br>011   | 11<br>3          | <ul> <li>True \$</li> <li>Sum \$</li> </ul> | Sum<br>with same bits |

| Unsigned | Binary |
|----------|--------|
| 0        | 000    |
| 1        | 001    |
| 2        | 010    |
| 3        | 011    |
| 4        | 100    |
| 5        | 101    |
| 6        | 110    |
| 7        | 111    |

## **Unsigned Addition in C**



| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

 Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

 Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)

|        |    | 010 | 2     |
|--------|----|-----|-------|
| Normal | +) | 101 | +) -3 |
| Case   |    | 111 | -1    |

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

| NI       | 010    | 2     |
|----------|--------|-------|
| Normal   | +) 101 | +) -3 |
| Case     | 111    | -1    |
| 0 1      | 110    | -2    |
| Overflow | +) 101 | +) -3 |
| Case     | 1011   | -5    |

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

|          | 010    | 2     |
|----------|--------|-------|
| Normal   | +) 101 | +) -3 |
| Case     | 111    | -1    |
| O        | 110    | -2    |
| Overtiow | +) 101 | +) -3 |
| Case     | 1011   | -5    |
|          | 011    | 3     |

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

|          | 010    | 2     |
|----------|--------|-------|
| Normal   | +) 101 | +) -3 |
| Case     | 111    | -1    |
| Overflow | 110    | -2    |
| Overnow  | +) 101 | +) -3 |
| Case     | 1011   | -5    |
|          | 011    | 3     |

| IS   | Signed | Binary |
|------|--------|--------|
|      | 0      | 000    |
| iye  | 1      | 001    |
|      | 2      | 010    |
|      | 3      | 011    |
| Min> | -4     | 100    |
|      | -3     | 101    |
|      | -2     | 110    |
|      | -1     | 111    |

Negative Overflow

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

| Normal           | 010<br>+) 101         | 2<br>+) -3       | Min            | 100<br>101<br>110 |
|------------------|-----------------------|------------------|----------------|-------------------|
| Case             | 111<br>110            | -1<br>-2         | -1             | 111<br><b>3</b>   |
| Overflow<br>Case | +) 101<br>1011<br>011 | +) -3<br>-5<br>3 | +) 001<br>0100 | +) 1<br>4         |

Signed

0

2

3

**Binary** 

000

001

010

011

Negative Overflow

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

| Normal<br>Case   | 010<br>+) 101      | 2<br>+) -3  | Min           | 100<br>101<br>110 |
|------------------|--------------------|-------------|---------------|-------------------|
|                  | 111                | -1          | -1            | 111               |
| Overflow<br>Case | 110<br>+) 101      | -2<br>+) -3 | 011<br>+) 001 | 3<br>+) 1         |
|                  | <b>1011</b><br>011 | -5<br>3     | 0100          | <b>4</b><br>- 4   |

Signed

0

2

3

**Binary** 

000

001

010

011

Negative Overflow

- Has identical bit-level behavior as unsigned addition (a big advantage over sign-magnitude)
- Overflow can also occur

| Normal<br>Case   |        |       | Min> -4 | 100  |
|------------------|--------|-------|---------|------|
|                  | 010    | 2     | -3      | 101  |
|                  | +) 101 | +) -3 | -2      | 110  |
|                  | 111    | -1    | -1      | 111  |
|                  |        |       |         |      |
| Overflow<br>Case | 110    | -2    | 011     | 3    |
|                  | +) 101 | +) -3 | +) 001  | +) 1 |
|                  | 1011   | -5    | 0100    | 4    |
|                  | 011    | 3     | 100     | -4   |
|                  |        |       |         |      |

**Negative Overflow** 

**Positive Overflow** 

Signed

U

2

З

Max

**Binary** 

000

001

010

011

Operands: *W* bits

True Sum: W+1 bits

Discard Carry: W bits



| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |

+) 111 +) 110 1101

|      | 111 |   |
|------|-----|---|
| +)   | 110 | _ |
| 1101 |     |   |

| Signed | Binary |
|--------|--------|
| 0      | 000    |
| 1      | 001    |
| 2      | 010    |
| 3      | 011    |
| -4     | 100    |
| -3     | 101    |
| -2     | 110    |
| -1     | 111    |



Cianad

|        |          |       | Signea | Binary |  |
|--------|----------|-------|--------|--------|--|
|        |          |       | 0      | 000    |  |
|        |          |       | 1      | 001    |  |
| 111    |          | -1    | 2      | 010    |  |
| +) 110 |          | +) -2 | 3      | 011    |  |
|        |          | ·// _ | -4     | 100    |  |
| 1101   |          | -3    | -3     | 101    |  |
|        | Truncate |       | -2     | 110    |  |
|        |          |       | -1     | 111    |  |



• This is not an overflow by definition



- This is not an overflow by definition
- Because the actual result can be represented using the bit width of the datatype (3 bits here)

• Goal: Computing Product of *w*-bit numbers *x*, *y* 

#### **Original Number (w bits)**

$$\begin{array}{c} OMax \ 2^{w-1}-1 \\ 0 \\ OMin \ -2^{w-1} \end{array}$$












## **Multiplication**

- Goal: Computing Product of *w*-bit numbers *x*, *y*
- Exact results can be bigger than w bits
  - Up to 2w bits (both signed and unsigned)

**Product** (2w bits) **Original Number (w bits) PMax** OMin<sup>2</sup> OMax <sup>2w -1</sup>– ، ا 0 OMin -2<sup>2w-2</sup> + 2<sup>w-1</sup> **PMin** OMin \* OMax

## **Unsigned Multiplication in C**



- Standard Multiplication Function
  - Ignores high order w bits
- Effectively Implements the following:

 $UMult_w(u, v) = u \cdot v \mod 2^w$ 

## **Unsigned Multiplication in C**



- Standard Multiplication Function
  - Ignores high order w bits

#### • Effectively Implements the following:

 $\mathsf{UMult}_w(u, v) = u \cdot v \mod 2^w$ 

|      |      | 1110 | 1001 |   | E9  |   | 223   |
|------|------|------|------|---|-----|---|-------|
| *    |      | 1101 | 0101 | * | D5  | * | 213   |
| **** | **** | 1101 | 1101 | C | 1DD | 4 | 47499 |
|      |      | 1101 | 1101 |   | DD  |   | 221   |

## Signed Multiplication in C



- Standard Multiplication Function
  - Ignores high order w bits
  - Some of which are different for signed vs. unsigned multiplication
  - Lower bits are the same

## Signed Multiplication in C



- Standard Multiplication Function
  - Ignores high order w bits
  - Some of which are different for signed vs. unsigned multiplication
  - Lower bits are the same

|           | 1110 | 1001 | E9   |   | -23 |
|-----------|------|------|------|---|-----|
| *         | 1101 | 0101 | * D5 | * | -43 |
| **** **** | 1101 | 1101 | 03DD |   | 989 |
|           | 1101 | 1101 | DD   |   | -35 |

- Operation
  - u << k gives u \*  $2^k$
  - Both signed and unsigned
     Operands: w bits



 $u \cdot 2^k$ 

k

. . .

0

00

0

0

...

010

U

 $2^k$ 

0

. . .

\*

- Operation
  - u << k gives u \*  $2^k$
  - Both signed and unsigned
     Operands: w bits

True Product: *w+k* bits

- Operation
  - $\mathbf{u} \ll \mathbf{k}$  gives  $\mathbf{u} \ast 2^k$
  - Both signed and unsig Operands: w bits

True Product: *w*+*k* bits

Discard k bits: w bits

| `             |   |       |     |     | k     |     |       |
|---------------|---|-------|-----|-----|-------|-----|-------|
| gned          |   | U     |     |     | • • • |     |       |
|               | * | $2^k$ | 0   | ••• | 010   | ••• | • 0 0 |
| $u \cdot 2^k$ |   | •     | ••• |     | 0     | ••• | • 0 0 |
|               |   |       |     | ••  |       | ••• | • 00  |

- Operation
  - u << k gives u \*  $2^k$
  - Both signed and unsigned
     Operands: w bits

True Product: <u>w+k</u> bits

Discard k bits: w bits

••• 0 0

n

00

0

k

0110

U

 $2^k$ 

\*

#### • Examples

- u << 3 == u \* 8
- $(u \ll 5) (u \ll 3) == u \ast 24$
- Most machines shift and add faster than multiply

 $u \cdot 2^k$ 

• Compiler generates this code automatically

### **Today: Representing Information in Binary**

- Why Binary (bits)?
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
  - Summary

## **Arithmetic: Basic Rules**

- Addition:
  - Unsigned/signed: Normal addition followed by truncate, same operation on bit level
- Multiplication:
  - Unsigned/signed: Normal multiplication followed by truncate, same operation on bit level
  - Shift: Power-of-2 Multiply

### Why Should I Use Unsigned?

- Don't use without understanding implications
  - Easy to make mistakes

```
unsigned int i;
for (i = cnt-2; i >= 0; i--)
a[i] += a[i+1];
```

• Can be very subtle

```
#define DELTA sizeof(int)
int i;
for (i = CNT; i-DELTA >= 0; i-= DELTA)
. . .
```

### Why Should I Use Unsigned? – Bit Set

• Use bits to represent my availability of the week

| b <sub>6</sub> | b <sub>5</sub> | b <sub>4</sub> | b <sub>3</sub> | b <sub>2</sub> | b <sub>1</sub> | b <sub>0</sub> |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| Sun            | Mon            | Tue            | Wed            | Thu            | Fri            | Sat            |
| 1              | 0              | 1              | 1              | 0              | 0              | 1              |

- Use 1 bit per day, 7 bits in total.
- If bit x is set to 1 then I'm available on day mapped to bit x.

### Why Should I Use Unsigned? – Bit Set

• Use bits to represent my availability of the week

| b <sub>6</sub> | $b_5$ | b <sub>4</sub> | b <sub>3</sub> | b <sub>2</sub> | b <sub>1</sub> | b <sub>0</sub> |
|----------------|-------|----------------|----------------|----------------|----------------|----------------|
| Sun            | Mon   | Tue            | Wed            | Thu            | Fri            | Sat            |
| 1              | 0     | 1              | 1              | 0              | 0              | 1              |

- Use 1 bit per day, 7 bits in total.
- If bit x is set to 1 then I'm available on day mapped to bit x.
- In C: unsigned int aval;

### Why Should I Use Unsigned? – Bit Set

• Use bits to represent my availability of the week

| b <sub>6</sub> | b <sub>5</sub> | b <sub>4</sub> | b <sub>3</sub> | b <sub>2</sub> | b <sub>1</sub> | b <sub>0</sub> |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| Sun            | Mon            | Tue            | Wed            | Thu            | Fri            | Sat            |
| 1              | 0              | 1              | 1              | 0              | 0              | 1              |

- Use 1 bit per day, 7 bits in total.
- If bit x is set to 1 then I'm available on day mapped to bit x.
- In C: unsigned int aval;

 $aval = 1^{*}2^{0} + 0^{*}2^{1} + 0^{*}2^{2} + 1^{*}2^{3} + 1^{*}2^{4} + 0^{*}2^{5} + 1^{*}2^{6} = 89_{10}$ 

## **Today: Floating Point**

- Background: Fractional binary numbers and fixed-point
- Floating point representation
- IEEE 754 standard
- Rounding, addition, multiplication
- Floating point in C
- Summary

- What does 10.012 mean?
  - C.f., Decimal

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

## $10.01_2 = 1^*2^1 + 0^*2^0 + 0^*2^{-1} + 1^*2^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

- What does 10.012 mean?
  - C.f., Decimal

## $12.45 = 1^*10^1 + 2^*10^0 + 4^*10^{-1} + 5^*10^{-2}$

# $10.01_2 = 1^*2^1 + 0^*2^0 + 0^*2^{-1} + 1^*2^{-2}$ $= 2.25_{10}$

## **Fractional Binary Numbers**



- What does 10.012 mean?
- C.f., Decimal
  - $12.45 = 1*10^1 + 2*10^0 + 4*10^{-1} + 5*10^{-2}$
- $10.01_2 = 1^*2^1 + 0^*2^0 + 0^*2^{-1} + 1^*2^{-2} = 2.25_{10}$

- What does 10.012 mean?
- C.f., Decimal
  - $12.45 = 1*10^1 + 2*10^0 + 4*10^{-1} + 5*10^{-2}$
- $10.01_2 = 1^*2^1 + 0^*2^0 + 0^*2^{-1} + 1^*2^{-2} = 2.25_{10}$



| Decimal | Binary |
|---------|--------|
| 0       | 0000   |
| 1       | 0001   |
| 2       | 0010   |
| 3       | 0011   |
| 4       | 0100   |
| 5       | 0101   |
| 6       | 0110   |
| 7       | 0111   |
| 8       | 1000   |
| 9       | 1001   |
| 10      | 1010   |
| 11      | 1011   |
| 12      | 1100   |
| 13      | 1101   |
| 14      | 1110   |
| 15      | 1111   |

- What does 10.012 mean?
- C.f., Decimal
  - $12.45 = 1*10^1 + 2*10^0 + 4*10^{-1} + 5*10^{-2}$
- $10.01_2 = 1^*2^1 + 0^*2^0 + 0^*2^{-1} + 1^*2^{-2} = 2.25_{10}$

#### 

| Decimal | Binary |
|---------|--------|
| 0       | 00.00  |
| 0.25    | 00.01  |
| 0.5     | 00.10  |
| 0.75    | 00.11  |
| 1       | 01.00  |
| 1.25    | 01.01  |
| 1.5     | 01.10  |
| 1.75    | 01.11  |
| 2       | 10.00  |
| 2.25    | 10.01  |
| 2.5     | 10.10  |
| 2.75    | 10.11  |
| 3       | 11.00  |
| 3.25    | 11.01  |
| 3.5     | 11.10  |
| 3.75    | 11.11  |

- What does 10.012 mean?
- C.f., Decimal
  - $12.45 = 1 \times 10^{1} + 2 \times 10^{0} + 4 \times 10^{-1} + 5 \times 10^{-2}$
- $10.01_2 = 1^{*}2^{1} + 0^{*}2^{0} + 0^{*}2^{-1} + 1^{*}2^{-2} = 2.25_{10}$

-----

$$\begin{array}{cccc}
01.10 & 1.50 \\
01.01 & + 1.25 \\
10.11 & 2.75 \end{array}$$

| Decimal | Binary |
|---------|--------|
| 0       | 00.00  |
| 0.25    | 00.01  |
| 0.5     | 00.10  |
| 0.75    | 00.11  |
| 1       | 01.00  |
| 1.25    | 01.01  |
| 1.5     | 01.10  |
| 1.75    | 01.11  |
| 2       | 10.00  |
| 2.25    | 10.01  |
| 2.5     | 10.10  |
| 2.75    | 10.11  |
| 3       | 11.00  |
| 3.25    | 11.01  |
| 3.5     | 11.10  |
| 3.75    | 11.11  |

- What does 10.012 mean?
- C.f., Decimal
  - $12.45 = 1*10^1 + 2*10^0 + 4*10^{-1} + 5*10^{-2}$
- $10.01_2 = 1^{*}2^{1} + 0^{*}2^{0} + 0^{*}2^{-1} + 1^{*}2^{-2} = 2.25_{10}$

#### <del>╋╪╪╪╋┊╡╝╝</del>

0 1 2 3

Integer Arithmetic Still Works!

| 01.10   | 1.50   |
|---------|--------|
| + 01.01 | + 1.25 |
| 10.11   | 2.75   |

| Decimal | Binary |
|---------|--------|
| 0       | 00.00  |
| 0.25    | 00.01  |
| 0.5     | 00.10  |
| 0.75    | 00.11  |
| 1       | 01.00  |
| 1.25    | 01.01  |
| 1.5     | 01.10  |
| 1.75    | 01.11  |
| 2       | 10.00  |
| 2.25    | 10.01  |
| 2.5     | 10.10  |
| 2.75    | 10.11  |
| 3       | 11.00  |
| 3.25    | 11.01  |
| 3.5     | 11.10  |
| 3.75    | 11.11  |
|         |        |

## **Fixed-Point Representation**

|  | Decimai | Binary |
|--|---------|--------|
| <ul> <li>Binary point stays fixed</li> </ul>                   | 0       | 00.00  |
| <ul> <li>Fixed interval between two representable</li> </ul>   | 0.25    | 00.01  |
| numbers as long as the binary point stays fixed                | 0.5     | 00.10  |
| <ul> <li>The interval in this example is 0.2510</li> </ul>     | 0.75    | 00.11  |
|  |         | 01.00  |
| Fixed-point representation of numbers                          | 1.25    | 01.01  |
| <ul> <li>Integer is one special case of fixed-point</li> </ul> | 1.5     | 01.10  |
|  | 1.75    | 01.11  |
|  | 2       | 10.00  |
| 0 1 2 3  | 2.25    | 10.01  |
|  | 2.5     | 10.10  |
|  | 2.75    | 10.11  |
|  | 3       | 11.00  |
|  | 3.25    | 11.01  |
|  | 3.5     | 11.10  |
|  | 3.75    | 11.11  |

## **Fixed-Point Representation**

Decimal Binary Binary point stays fixed 0000.U Fixed interval between two representable 0001. numbers as long as the binary point stays fixed 2 0010. 3 0011. The interval in this example is 0.2510 4 0100. Fixed-point representation of numbers 0101. 5 0110. Integer is one special case of fixed-point 6 0111. 7 8 1000. 7 Ω 2 3 Δ 5 6 15 9 1001. 10 1010. 11 1011. 12 1100. 13 1101. 14 1110.

1111.

15
• Can exactly represent numbers only of the form x/2<sup>k</sup>

• Can exactly represent numbers only of the form x/2<sup>k</sup>



- Can exactly represent numbers only of the form x/2<sup>k</sup>
  - Other rational numbers have repeating bit representations



- Can exactly represent numbers only of the form x/2<sup>k</sup>
  - Other rational numbers have repeating bit representations

| Decimal Value | <b>Binary Representation</b> |
|---------------|------------------------------|
| 1/3           | 0.0101010101[01]             |
| 1/5           | 0.001100110011[0011]         |
| 1/10          | 0.0001100110011[0011]        |



 Can't represent very small and very large numbers at the same time

- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers

- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers



- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers



- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers



- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers
  - To represent very small numbers, the (fixed) interval needs to be small, making it hard to represent large numbers



- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers
  - To represent very small numbers, the (fixed) interval needs to be small, making it hard to represent large numbers



- Can't represent very small and very large numbers at the same time
  - To represent very large numbers, the (fixed) interval needs to be large, making it hard to represent small numbers
  - To represent very small numbers, the (fixed) interval needs to be small, making it hard to represent large numbers

Unrepresentable large numbers ++∞ 0 A Small Number

# **Today: Floating Point**

- Background: Fractional binary numbers and fixed-point
- Floating point representation
- IEEE 754 standard
- Rounding, addition, multiplication
- Floating point in C
- Summary

- In decimal:  $M \times 10^{E}$ 
  - E is an integer
  - Normalized form:  $1 \le |M| \le 10$

- In decimal:  $M \times 10^{E}$ 
  - E is an integer
  - Normalized form:  $1 \le |M| \le 10$

| Decimal Value    | Scientific Notation       |
|------------------|---------------------------|
| 2                | 2×10 <sup>0</sup>         |
| -4,321.768       | -4.321768×10 <sup>3</sup> |
| 0.000 000 007 51 | 7.51×10 <sup>-9</sup>     |

- In decimal:  $M \times 10^{E}$ 
  - E is an integer
  - Normalized form:  $1 \le |M| \le 10$

# $M \times 10^{E}$

| Decimal Value    | Scientific Notation       |
|------------------|---------------------------|
| 2                | 2×10 <sup>0</sup>         |
| -4,321.768       | -4.321768×10 <sup>3</sup> |
| 0.000 000 007 51 | 7.51×10 <sup>-9</sup>     |

- In decimal: M × 10<sup>E</sup>
  - E is an integer
  - Normalized form:  $1 \le |M| \le 10$



#### Significand

| Decimal Value    | Scientific Notation       |
|------------------|---------------------------|
| 2                | 2×10 <sup>0</sup>         |
| -4,321.768       | -4.321768×10 <sup>3</sup> |
| 0.000 000 007 51 | 7.51×10 <sup>-9</sup>     |

- In decimal:  $M \times 10^{E}$ 
  - *E* is an integer
  - Normalized form:  $1 \le |M| \le 10$



Significand Base

| Decimal Value    | Scientific Notation       |
|------------------|---------------------------|
| 2                | 2×10 <sup>0</sup>         |
| -4,321.768       | -4.321768×10 <sup>3</sup> |
| 0.000 000 007 51 | 7.51×10 <sup>-9</sup>     |

- In decimal:  $M \times 10^{E}$ 
  - *E* is an integer
  - Normalized form:  $1 \le |M| \le 10$



Significand Base

| Decimal Value    | Scientific Notation       |
|------------------|---------------------------|
| 2                | 2×10 <sup>0</sup>         |
| -4,321.768       | -4.321768×10 <sup>3</sup> |
| 0.000 000 007 51 | 7.51×10 <sup>-9</sup>     |

- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction



| Binary Value  | Scientific Notation                                |
|---------------|--|
| 1110110110110 | (-1) <sup>0</sup> 1.110110110110 x 2 <sup>12</sup> |
| -101.11       | (-1) <sup>1</sup> 1.0111 x 2 <sup>2</sup>          |
| 0.00101       | (-1) <sup>0</sup> 1.01 x 2 <sup>-3</sup>           |

- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction



- If I tell you that there is a number where:
  - Fraction = 0101
  - s = 1
  - E = 10
  - You could reconstruct the number as (-1)<sup>1</sup>1.0101x2<sup>10</sup>

- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction



- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction
- Encoding



- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction
- Encoding



- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction
- Encoding
  - MSB **s** is sign bit <mark>s</mark>





- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction
- Encoding
  - MSB **s** is sign bit s
  - exp field encodes Exponent (but not exactly the same, more later)





- In binary: (-1)<sup>s</sup> M 2<sup>E</sup>
- Normalized form:
  - 1<=*M* < 2
  - $M = 1.bob_1b_2b_3...$ Fraction
- Encoding
  - MSB **s** is sign bit s
  - *exp* field encodes Exponent (but not exactly the same, more later)
  - frac field encodes Fraction (but not exactly the same, more later)









• *exp* has 3 bits, interpreted as an unsigned value



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):


- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):
  - bias = 3



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):
  - bias = 3
  - If *E* = -2, *exp* is 1 (001<sub>2</sub>)



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):
  - bias = 3
  - If *E* = -2, *exp* is 1 (001<sub>2</sub>)

| E  | ехр |
|----|-----|
| -3 | 000 |
| -2 | 001 |
| -1 | 010 |
| С  | 011 |
| 1  | 100 |
| 2  | 101 |
| 3  | 110 |
| 4  | 111 |



- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):
  - bias = 3
  - If *E* = -2, *exp* is 1 (001<sub>2</sub>)
  - Reserve 000 and 111 for other purposes (more on this later)





- *exp* has 3 bits, interpreted as an unsigned value
  - If exp were E, we could represent exponents from 0 to 7
  - How about negative exponent?
  - Subtract a bias term: *E* = *exp bias* (i.e., exp = E + bias)
  - bias is always 2<sup>k-1</sup> 1, where k is number of exponent bits
- Example when we use 3 bits for exp (i.e., k = 3):
  - bias = 3
  - If *E* = -2, *exp* is 1 (001<sub>2</sub>)
  - Reserve 000 and 111 for other purposes (more on this later)
  - We can now represent exponents from -2 (exp 001) to 3 (exp 110)





- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10



- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

# $-10.1_2 = (-1)^1 \ 1.01 \ x \ 2^1$



- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^1 1.01 \times 2^1$$



- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^1 1.01 \times 2^1$$

Į



- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^1 1.01 \times 2^1$$



- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^{1} 1.01 \times 2^{1}$$





- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^{1} 1.01 \times 2^{1}$$





- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^{1} 1.01 \times 2^{1}$$

I





- frac has 2 bits, append them after "1." to form M
  - *frac* = 10 implies M = 1.10
- Putting it Together: An Example:

$$-10.1_2 = (-1)^{1} 1.01 \times 2^{1}$$

I

