

**Princeton University**  
Computer Science 217: Introduction to Programming Systems



## Number Systems and Number Representation

Q: Why do computer programmers confuse Christmas and Halloween?  
A: Because 25 Dec = 31 Oct



1

1

**Goals of this Lecture**



Help you learn (or refresh your memory) about:

- The binary, hexadecimal, and octal number systems
- Finite representation of unsigned integers
- Finite representation of signed integers
- Finite representation of rational (floating-point) numbers

Why?

- A power programmer must know number systems and data representation to fully understand C's **primitive data types**

/ Primitive values and the operations on them

2

2

**Agenda**



### Number Systems

Finite representation of unsigned integers  
Finite representation of signed integers  
Finite representation of rational (floating-point) numbers

3

3

**The Decimal Number System**



**Name**  
• "decem" (Latin)  $\Rightarrow$  ten

**Characteristics**

- Ten symbols
  - 0 1 2 3 4 5 6 7 8 9
- Positional
  - $2945 \neq 2495$
  - $2945 = (2*10^3) + (9*10^2) + (4*10^1) + (5*10^0)$

(Most) people use the decimal number system



4

4

**The Binary Number System**



**binary**  
*adjective:* being in a state of one of two mutually exclusive conditions such as on or off, true or false, molten or frozen, presence or absence of a signal.  
From Late Latin *biniārius* ("consisting of two").

**Characteristics**

- Two symbols: 0 1
- Positional:  $1010_b \neq 1100_b$

Most (digital) computers use the binary number system

**Terminology**

- Bit:** a binary digit
- Byte:** (typically) 8 bits
- Nibble (or nybble):** 4 bits



5

5

**Decimal-Binary Equivalence**



Decimal	Binary	Decimal	Binary
0	0	16	10000
1	1	17	10001
2	10	18	10010
3	11	19	10011
4	100	20	10100
5	101	21	10101
6	110	22	10110
7	111	23	10111
8	1000	24	11000
9	1001	25	11001
10	1010	26	11010
11	1011	27	11011
12	1100	28	11100
13	1101	29	11101
14	1110	30	11110
15	1111	31	11111
		...	...

6

1

## Decimal-Binary Conversion

Binary to decimal: expand using positional notation

$$\begin{aligned}100101_2 &= (1 \cdot 2^5) + (0 \cdot 2^4) + (0 \cdot 2^3) + (1 \cdot 2^2) + (0 \cdot 2^1) + (1 \cdot 2^0) \\&= 32 + 0 + 0 + 4 + 0 + 1 \\&= 37\end{aligned}$$

Most-significant bit (msb)

Least-significant bit (lsb)

7

## Integer Decimal-Binary Conversion

Integer

Binary to decimal: expand using positional notation

$$\begin{aligned}100101_2 &= (1 \cdot 2^5) + (0 \cdot 2^4) + (0 \cdot 2^3) + (1 \cdot 2^2) + (0 \cdot 2^1) + (1 \cdot 2^0) \\&= 32 + 0 + 0 + 4 + 0 + 1 \\&= 37\end{aligned}$$

These are integers

They exist as their pure selves  
no matter how we might choose  
to represent them with our  
fingers or toes

8

## Integer-Binary Conversion

Integer to binary: do the reverse

- Determine largest power of 2 that's  $\leq$  number; write template

$$37 = (? \cdot 2^5) + (? \cdot 2^4) + (? \cdot 2^3) + (? \cdot 2^2) + (? \cdot 2^1) + (? \cdot 2^0)$$

- Fill in template

$$\begin{array}{r} 37 = (1 \cdot 2^5) + (0 \cdot 2^4) + (0 \cdot 2^3) + (1 \cdot 2^2) + (0 \cdot 2^1) + (1 \cdot 2^0) \\ -32 \\ \hline 5 \\ -4 \\ \hline 1 \\ -1 \\ \hline 0 \end{array} \quad 100101_2$$

9

## Integer-Binary Conversion

Integer to binary shortcut

- Repeatedly divide by 2, consider remainder

$$\begin{array}{r} 37 \div 2 = 18 \text{ R } 1 \\ 18 \div 2 = 9 \text{ R } 0 \\ 9 \div 2 = 4 \text{ R } 1 \\ 4 \div 2 = 2 \text{ R } 0 \\ 2 \div 2 = 1 \text{ R } 0 \\ 1 \div 2 = 0 \text{ R } 1 \end{array}$$

Read from bottom  
to top:  $100101_2$

10

## The Hexadecimal Number System

### Name

- "hexa-" (Ancient Greek ἑξα-)  $\Rightarrow$  six
- "decem" (Latin)  $\Rightarrow$  ten

### Characteristics

- Sixteen symbols
  - 0 1 2 3 4 5 6 7 8 9 A B C D E F
- Positional
- $A13D_{10} \neq 3DA1_8$

Computer programmers often use hexadecimal or "hex"

- In C: `0x` prefix (`0xA13D`, etc.)

Why?

11

## Decimal-Hexadecimal Equivalence

Decimal	Hex
0	0
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	A
11	B
12	C
13	D
14	E
15	F

Decimal	Hex
16	10
17	11
18	12
19	13
20	14
21	15
22	16
23	17
24	18
25	19
26	1A
27	1B
28	1C
29	1D
30	1E
31	1F

Decimal	Hex
32	20
33	21
34	22
35	23
36	24
37	25
38	26
39	27
40	28
41	29
42	2A
43	2B
44	2C
45	2D
46	2E
47	2F
...	...

12

## Integer-Hexadecimal Conversion

Hexadecimal to integer: expand using positional notation

$$\begin{aligned} 25_{\text{H}} &= (2 * 16^1) + (5 * 16^0) \\ &= 32 + 5 \\ &= 37 \end{aligned}$$

Integer to hexadecimal: use the shortcut

$$\begin{array}{r} 37 \quad / \quad 16 = 2 \text{ R } 5 \\ 2 \quad / \quad 16 = 0 \text{ R } 2 \end{array}$$

Read from bottom to top:  $25_{\text{H}}$



13

## Binary-Hexadecimal Conversion

Observation:  $16^1 = 2^4$

- Every 1 hexadecimal digit corresponds to 4 binary digits

Binary to hexadecimal

$$\begin{array}{cccc} 1010000100111101_2 \\ \text{A} \quad 1 \quad 3 \quad \text{D}_H \end{array}$$

Digit count in binary number  
not a multiple of 4  $\Rightarrow$   
pad with zeros on left

Hexadecimal to binary

$$\begin{array}{cccc} \text{A} \quad 1 \quad 3 \quad \text{D}_H \\ 1010000100111101_2 \end{array}$$

Discard leading zeros from  
binary number if appropriate

Is it clear why programmers  
often use hexadecimal?

14

## iClicker Question

Q: Convert binary 101010 into decimal and hex

- A. 21 decimal, 1A hex
- B. 42 decimal, 2A hex
- C. 48 decimal, 32 hex
- D. 55 decimal, 4G hex

Hint: convert to hex first

15

## The Octal Number System



Name

- "octo" (Latin)  $\Rightarrow$  eight

Characteristics

- Eight symbols
  - 0 1 2 3 4 5 6 7
- Positional
  - $1743_8 \neq 7314_8$

Computer programmers often use octal (so does Mickey!)

- In C: 0 prefix (01743, etc.)

Why?

16

## Agenda



Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational (floating-point) numbers

17

17

## Integral Types in Java vs. C



	Java	C
Unsigned types	char // 16 bits	unsigned char /* Note 2 */ unsigned short unsigned (int) unsigned long
Signed types	byte // 8 bits short // 16 bits int // 32 bits long // 64 bits	signed char /* Note 2 */ (signed) short (signed) int (signed) long

1. Not guaranteed by C, but on arm lab, char = 8 bits, short = 16 bits, int = 32 bits, long = 64 bits

2. Not guaranteed by C, but on arm lab, char is unsigned

To understand C, must consider representation of both unsigned and signed integers

18

18

## Representing Unsigned Integers

### Mathematics

- Range is 0 to  $\infty$

### Computer programming

- Range limited by computer's **word size**
- Word size is n bits  $\Rightarrow$  range is 0 to  $2^n - 1$
- Exceed range  $\Rightarrow$  **overflow**

### Typical computers today

- $n = 32$  or  $64$ , so range is 0 to  $2^{32} - 1$  or  $2^{64} - 1$  (huge!)

### Pretend computer

- $n = 4$ , so range is 0 to  $2^4 - 1$  (15)

### Hereafter, assume word size = 4

- All points generalize to word size = 64, word size = n



19

## Representing Unsigned Integers

### On pretend computer

Unsigned Integer	Rep
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



20

## Adding Unsigned Integers

### Addition

$$\begin{array}{r} & 1 \\ 3 & + 10 \\ + & 10 \\ \hline 13 & 1101_2 \end{array}$$

Start at right column  
Proceed leftward  
Carry 1 when necessary

$$\begin{array}{r} & 111 \\ 7 & + 10 \\ + & 10 \\ \hline 1 & 0001_2 \end{array}$$

Beware of overflow

How would you  
detect overflow  
programmatically?

Results are mod  $2^4$ 

21



## Subtracting Unsigned Integers

### Subtraction

$$\begin{array}{r} & 111 \\ 10 & - 7 \\ - & 7 \\ \hline 3 & 0011_2 \end{array}$$

Start at right column  
Proceed leftward  
Borrow when necessary

$$\begin{array}{r} & 1 \\ 3 & - 10 \\ - & 10 \\ \hline 9 & 1001_2 \end{array}$$

Beware of overflow

How would you  
detect overflow  
programmatically?

Results are mod  $2^4$ 

22



## Shifting Unsigned Integers

### Bitwise right shift (`>>` in C): fill on left with zeros

$$\begin{array}{l} 10 \gg 1 \Rightarrow 5 \\ 1010_2 \quad 0101_2 \end{array}$$

What is the effect  
arithmetically?

$$\begin{array}{l} 10 \gg 2 \Rightarrow 2 \\ 1010_2 \quad 0010_2 \end{array}$$

### Bitwise left shift (`<<` in C): fill on right with zeros

$$\begin{array}{l} 5 \ll 1 \Rightarrow 10 \\ 0101_2 \quad 1010_2 \end{array}$$

What is the effect  
arithmetically?

$$\begin{array}{l} 3 \ll 2 \Rightarrow 12 \\ 0011_2 \quad 1100_2 \end{array}$$

Results are mod  $2^4$ 

23



### Bitwise NOT (~ in C)

- Flip each bit

$$\begin{array}{l} -10 \Rightarrow 5 \\ 1010_2 \quad 0101_2 \end{array}$$

$$\begin{array}{l} -5 \Rightarrow 10 \\ 0101_2 \quad 1010_2 \end{array}$$

### Bitwise AND (& in C)

- Logical AND corresponding bits

$$\begin{array}{r} 10 \\ \& 7 \\ - & 0111_2 \\ \hline 2 & 0010_2 \end{array}$$

$$\begin{array}{r} 10 \\ \& 2 \\ - & 0010_2 \\ \hline 2 & 0010_2 \end{array}$$

Useful for "masking" bits to 0

24

23

24

## Other Operations on Unsigned Ints



### Bitwise OR: (| in C)

- Logical OR corresponding bits

10	1010 <sub>B</sub>
1	0001 <sub>B</sub>
--	----
11	1011 <sub>B</sub>

Useful for "masking" bits to 1

### Bitwise exclusive OR (^ in C)

- Logical exclusive OR corresponding bits

10	1010 <sub>B</sub>
^ 10	^ 1010 <sub>B</sub>
--	----
0	0000 <sub>B</sub>

x ^ x sets  
all bits to 0

25

25

## iClicker Question

Q: How do you set bit "n" (counting lsb=0) of **unsigned** variable "u" to zero?

- A. u &= (0 << n);
- B. u |= (1 << n);
- C. u &= ~(1 << n);
- D. u |= ~(1 << n);
- E. u = ~u ^ (1 << n);

26

## Aside: Using Bitwise Ops for Arith



Can use <<, >>, and & to do some arithmetic efficiently

$$x * 2^y == x \ll y$$

- $3 * 4 = 3 * 2^2 = 3 \ll 2 \Rightarrow 12$

Fast way to multiply  
by a power of 2

$$x / 2^y == x \gg y$$

- $13/4 = 13/2^2 = 13 \gg 2 \Rightarrow 3$

Fast way to divide  
unsigned by power of 2

$$x \% 2^y == x \& (2^y - 1)$$

- $13 \% 4 = 13 \% 2^2 = 13 \& (2^2 - 1)$
- $= 13 \& 3 \Rightarrow 1$

Fast way to mod  
by a power of 2

13	1101 <sub>B</sub>
& 3	& 0011 <sub>B</sub>
--	----
1	0001 <sub>B</sub>

Many compilers will  
do these transformations  
automatically!

27

27

## Aside: Example C Program



```
#include <stdio.h>
#include <stdlib.h>
int main(void)
{
    unsigned int n;
    unsigned int count = 0;
    printf("Enter an unsigned integer: ");
    if (scanf("%u", &n) != 1)
    {
        fprintf(stderr, "Error: Expect unsigned int.\n");
        exit(EXIT_FAILURE);
    }
    while (n > 0)
    {
        count += (n & 1);
        n = n >> 1;
    }
    printf("%u\n", count);
    return 0;
}
```

What does it  
write?

How could you  
express this more  
succinctly?

28

28

## Agenda



Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational (floating-point) numbers

29

29

## Sign-Magnitude



Integer	Rep
-7	1111
-6	1110
-5	1101
-4	1100
-3	1011
-2	1010
-1	1001
0	1000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

### Definition

High-order bit indicates sign

0 ⇒ positive

1 ⇒ negative

Remaining bits indicate magnitude

$0101_B = 101_B = 5$

$1101_B = -101_B = -5$

30

5

## Sign-Magnitude (cont.)



Integer	Rep
-7	1111
-6	1110
-5	1101
-4	1100
-3	1011
-2	1010
-1	1001
0	1000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

### Computing negative

$\text{neg}(x) = \text{flip high order bit of } x$

$$\text{neg}(0101_B) = 1101_B$$

$$\text{neg}(1101_B) = 0101_B$$

### Pros and cons

- + easy to understand, easy to negate
- + symmetric
- two representations of zero
- need different algorithms to add signed and unsigned numbers

31

## Ones' Complement



Integer	Rep
-7	1000
-6	1001
-5	1010
-4	1011
-3	1100
-2	1101
-1	1110
0	1111
1	0000
2	0001
3	0010
4	0011
5	0100
6	0101
7	0110

### Definition

High-order bit has weight  $-(2^{b-1})$

$$1010_B = (1 * -7) + (0 * 4) + (1 * 2) + (0 * 1) = -5$$

$$0010_B = (0 * -7) + (0 * 4) + (1 * 2) + (0 * 1) = 2$$

32

31

32

## Ones' Complement (cont.)



Integer	Rep
-7	1000
-6	1001
-5	1010
-4	1011
-3	1100
-2	1101
-1	1110
0	1111
1	0000
2	0001
3	0010
4	0011
5	0100
6	0101
7	0110

### Computing negative

$\text{neg}(x) = \sim x$

$$\text{neg}(0101_B) = 1010_B$$

$$\text{neg}(1010_B) = 0101_B$$

### Similar pros and cons to sign-magnitude

33

## Two's Complement



Integer	Rep
-8	1000
-7	1001
-6	1010
-5	1011
-4	1100
-3	1101
-2	1110
-1	1111
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

### Definition

High-order bit has weight  $-(2^{b-1})$

$$1010_B = (1 * -8) + (0 * 4) + (1 * 2) + (0 * 1) = -6$$

$$0010_B = (0 * -8) + (0 * 4) + (1 * 2) + (0 * 1) = 2$$

34

33

34

## Two's Complement (cont.)



Integer	Rep
-8	1000
-7	1001
-6	1010
-5	1011
-4	1100
-3	1101
-2	1110
-1	1111
0	0000
1	0001
2	0010
3	0011
4	0100
5	0101
6	0110
7	0111

### Computing negative

$\text{neg}(x) = \sim x + 1$

$$\text{neg}(x) = \text{onescomp}(x) + 1$$

$$\text{neg}(0101_B) = 1010_B + 1 = 1011_B$$

$$\text{neg}(1011_B) = 0100_B + 1 = 0101_B$$

### Pros and cons

- not symmetric ("extra" negative number)
- + one representation of zero
- + same algorithm adds unsigned numbers or signed numbers

35

35

## Two's Complement (cont.)



Almost all computers today use two's complement to represent signed integers

- Arithmetic is easy!

Is it after 1980?  
OK, then we're surely  
two's complement



Hereafter, assume two's complement

36

36

## Adding Signed Integers

pos + pos

$$\begin{array}{r} 3 \\ + 3 \\ \hline 6 \end{array}$$

0011<sub>b</sub>      0011<sub>b</sub>  
+    +  
-----  
0110<sub>b</sub>

pos + pos (overflow)

$$\begin{array}{r} 7 \\ + 1 \\ \hline -8 \end{array}$$

0111<sub>b</sub>      0001<sub>b</sub>  
+    +  
-----  
1000<sub>b</sub>

pos + neg

$$\begin{array}{r} 3 \\ + -1 \\ \hline 2 \end{array}$$

0011<sub>b</sub>      1111<sub>b</sub>  
+    +  
-----  
0010<sub>b</sub>

How would you detect overflow programmatically?

neg + neg

$$\begin{array}{r} -3 \\ + -2 \\ \hline -5 \end{array}$$

1101<sub>b</sub>      1110<sub>b</sub>  
+    +  
-----  
1011<sub>b</sub>

neg + neg (overflow)

$$\begin{array}{r} -6 \\ + -5 \\ \hline 5 \end{array}$$

1011<sub>b</sub>      1011<sub>b</sub>  
+    +  
-----  
0101<sub>b</sub>



37

## Subtracting Signed Integers

Perform subtraction with borrows

$$\begin{array}{r} 3 \\ - 4 \\ \hline -1 \end{array}$$

0011<sub>b</sub>      0100<sub>b</sub>  
-    -  
-----  
1111<sub>b</sub>

Compute two's comp and add

$$\begin{array}{r} 3 \\ + -4 \\ \hline -1 \end{array}$$

0011<sub>b</sub>      1100<sub>b</sub>  
+    +  
-----  
1111<sub>b</sub>

$$\begin{array}{r} -5 \\ - 2 \\ \hline -7 \end{array}$$

1011<sub>b</sub>      0010<sub>b</sub>  
-    -  
-----  
1001<sub>b</sub>

38

38

## Negating Signed Ints: Math



Question: Why does two's comp arithmetic work?

Answer:  $[-b] \bmod 2^4 = [\text{twoscomp}(b)] \bmod 2^4$ 

$$\begin{aligned} [-b] \bmod 2^4 &= [2^4 - b] \bmod 2^4 \\ &= [2^4 - 1 - b + 1] \bmod 2^4 \\ &= [(2^4 - 1 - b) + 1] \bmod 2^4 \\ &= [\text{onescomp}(b) + 1] \bmod 2^4 \\ &= [\text{twoscomp}(b)] \bmod 2^4 \end{aligned}$$

See Bryant &amp; O'Hallaron book for much more info

39

## Pithier Rationale: Math



Ring theory.

If  $n > 0$ ,  $\mathbb{Z}/(n)$  is a finite commutative ring, with properties:

$$\bar{a}_n + \bar{b}_n = \overline{(a+b)_n}; \bar{a}_n - \bar{b}_n = \overline{(a-b)_n}; \bar{a}_n \bar{b}_n = \overline{(ab)_n}$$

41

## Subtracting Signed Ints: Math



And so:

$$[a - b] \bmod 2^4 = [a + \text{twoscomp}(b)] \bmod 2^4$$

$$\begin{aligned} [a - b] \bmod 2^4 &= [a + 2^4 - b] \bmod 2^4 \\ &= [a + 2^4 - 1 - b + 1] \bmod 2^4 \\ &= [a + (2^4 - 1 - b) + 1] \bmod 2^4 \\ &= [a + \text{onescomp}(b) + 1] \bmod 2^4 \\ &= [a + \text{twoscomp}(b)] \bmod 2^4 \end{aligned}$$

See Bryant &amp; O'Hallaron book for much more info

40

## Shifting Signed Integers

Bitwise left shift ( $<<$  in C): fill on right with zeros

$$\begin{array}{r} 3 \ll 1 \Rightarrow 6 \\ 0011_8 \quad 0110_8 \\ -3 \ll 1 \Rightarrow -6 \\ 1101_8 \quad 1010_8 \end{array}$$

What is the effect arithmetically?

Results are mod  $2^4$ 

Bitwise right shift: fill on left with ???

42

42

41

## Shifting Signed Integers (cont.)

Bitwise arithmetic right shift: fill on left with sign bit

6 >> 1 => 3
0110 <sub>b</sub> 0011 <sub>b</sub>
-6 >> 1 => -3
1010 <sub>b</sub> 1101 <sub>b</sub>

What is the effect arithmetically?

Bitwise logical right shift: fill on left with zeros

6 >> 1 => 3
0110 <sub>b</sub> 0011 <sub>b</sub>
-6 >> 1 => 5
1010 <sub>b</sub> 0101 <sub>b</sub>

What is the effect arithmetically ???

In C, right shift (>>) could be logical or arithmetic

- Not specified by standard (happens to be arithmetic on armab)
- Best to avoid shifting signed integers**

43

## Other Operations on Signed Ints

Bitwise NOT (~ in C)

- Same as with unsigned ints

Bitwise AND (& in C)

- Same as with unsigned ints

Bitwise OR: (| in C)

- Same as with unsigned ints

Bitwise exclusive OR (^ in C)

- Same as with unsigned ints

**Best to avoid with signed integers**



44

43

44

## Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

**Finite representation of rational (floating-point) numbers**

45

45

## Rational Numbers



### Mathematics

- A **rational** number is one that can be expressed as the **ratio** of two integers
- Unbounded range and precision

### Computer science

- Finite range and precision
- Approximate using **floating point** number

46

46

## Floating Point Numbers



Like scientific notation: e.g., c is  
 $2.99792458 \times 10^8$  m/s

This has the form  
 $(\text{multiplier}) \times (\text{base})^{(\text{power})}$

In the computer,

- Multiplier** is called mantissa
- Base is almost always 2
- Power** is called exponent

47

## IEEE Floating Point Representation



Common finite representation: **IEEE floating point**

- More precisely: ISO/IEEE 754 standard

Using 32 bits (type **float** in C):

- 1 bit: sign (0=>positive, 1=>negative)
- 8 bits: exponent + 127
- 23 bits: binary fraction of the form 1.bbbbbbb...bbbbbbb...bbb

Using 64 bits (type **double** in C):

- 1 bit: sign (0=>positive, 1=>negative)
- 11 bits: exponent + 1023
- 52 bits: binary fraction of the form 1.bbbbbbb...bbbbbbb...bbb...bbb

48

48

## Floating Point Example



### Sign (1 bit):

- 1  $\Rightarrow$  negative

**32-bit representation**

### Exponent (8 bits):

- $10000011_2 = 131$
- $131 - 127 = 4$

### Mantissa (23 bits):

- $1.10110110000000000000000_2$
- $1 + (1 \cdot 2^{-1}) + (0 \cdot 2^{-2}) + (1 \cdot 2^{-3}) + (1 \cdot 2^{-4}) + (0 \cdot 2^{-5}) + (1 \cdot 2^{-6}) + (1 \cdot 2^{-7}) = 1.7109375$

### Number:

$$\bullet -1.7109375 \cdot 2^4 = -27.375$$

49

## When was floating-point invented?



Answer: long before computers!

### mantissa

*noun*  
decimal part of a logarithm, 1865, from Latin *mantissa* "a worthless addition, make-weight," perhaps a Gaulish word introduced into Latin via Etruscan (cf. Old Irish *meit*, Welsh *mant* "size").

### COMMON LOGARITHMS $\log_{10}x$

x	0	1	2	3	4	5	6	7	8	9	$\Delta_m$	I	2	3
50	-6990	6998	7007	7016	7024	7033	7042	7050	7059	7067	9	1	2	3
51	-7076	7084	7093	7101	7110	7118	7126	7135	7143	7152	8	1	2	2
52	-7160	7168	7177	7185	7193	7202	7210	7218	7226	7235	8	1	2	2
53	-7243	7251	7259	7267	7275	7284	7292	7300	7308	7316	8	1	2	2
54	-7324	7332	7340	7348	7356	7364	7372	7380	7388	7396	8	1	2	2
55	-7404	7412	7419	7427	7435	7443	7451	7459	7466	7474	8	1	2	2

50

## Floating Point Consequences



"Machine epsilon": smallest positive number you can add to 1.0 and get something other than 1.0

For float:  $\epsilon \approx 10^{-7}$

- No such number as 1.0000000001
- Rule of thumb: "almost 7 digits of precision"

For double:  $\epsilon \approx 2 \times 10^{-16}$

- Rule of thumb: "not quite 16 digits of precision"

These are all *relative* numbers

51

## Floating Point Consequences, cont



Just as decimal number system can represent only some rational numbers with finite digit count...

- Example:  $1/3$  **cannot** be represented

Binary number system can represent only some rational numbers with finite digit count

- Example:  $1/5$  **cannot** be represented

### Beware of roundoff error

- Error resulting from inexact representation
- Can accumulate
- Be careful when comparing two floating-point numbers for equality

Decimal	Rational
.3	$3/10$
.33	$33/100$
.333	$333/1000$
...	

Binary	Rational
0.0	$0/2$
0.01	$1/4$
0.010	$2/8$
0.0011	$3/16$
0.00110	$6/32$
0.001101	$13/64$
0.0011010	$26/128$
0.00110011	$51/256$
...	

52

## iClicker Question

Q: What does the following code print?

```
double sum = 0.0;
int i;
for (i = 0; i < 10; i++)
    sum += 0.1;
if (sum == 1.0)
    printf("All good!\n");
else
    printf("Yikes!\n");
```

- A. All good!
- B. Yikes!
- C. Code crashes
- D. Code enters an infinite loop

53

## Summary



The binary, hexadecimal, and octal number systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational (floating-point) numbers

### Essential for proper understanding of

- C primitive data types
- Assembly language
- Machine language

54