Princeton University

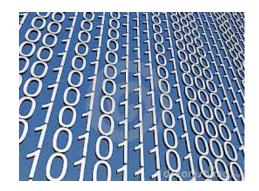
Computer Science 217: Introduction to Programming Systems



Number Systems and Number Representation

Q: Why do computer programmers confuse Christmas and Halloween?

A: Because 25 Dec = 31 Oct



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

Agenda



Number Systems

Finite representation of unsigned integers

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



The Decimal Number System

Name

• "decem" (Latin) \Rightarrow ten

Characteristics

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

(Most) people use the decimal number system



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 *bīnārius* ("consisting of two").

Characteristics

- Two symbols
 - 0 1
- Positional
 - 1010_B ≠ 1100_B

Most (digital) computers use the binary number system

Terminology

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

Why?



Decimal-Binary Equivalence

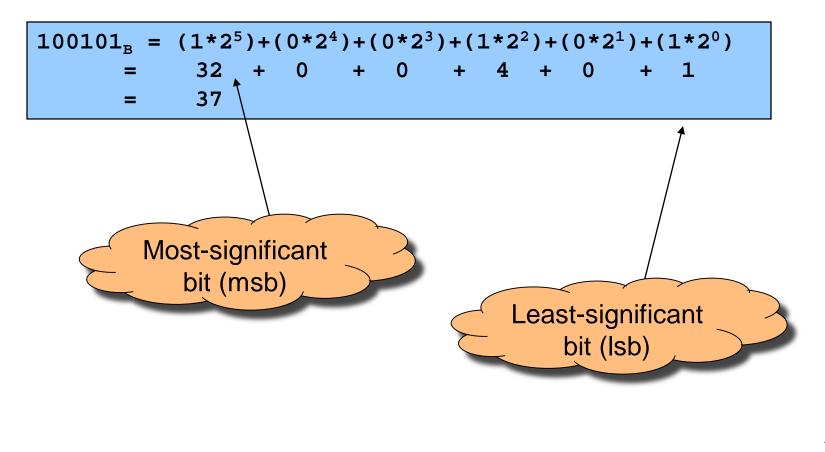
D

ecimal	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
		I	• • •	•••

Decimal-Binary Conversion



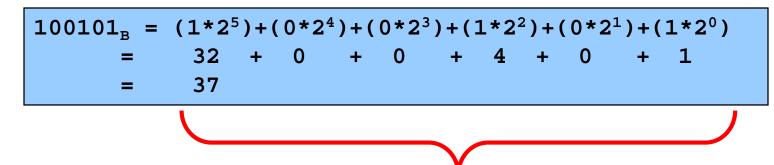
Binary to decimal: expand using positional notation



Integer Decimal-Binary Conversion



Integer Binary to decimal: expand using positional notation



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

Integer-Binary Conversion



Integer to binary: do the reverse

• Determine largest power of 2 that's ≤ number; write template

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

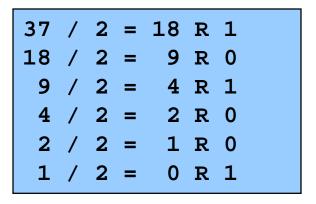
• Fill in template

 $37 = (1*2^{5})+(0*2^{4})+(0*2^{3})+(1*2^{2})+(0*2^{1})+(1*2^{0})$ -325 -41
100101_B
-1
0

Integer-Binary Conversion

Integer to binary shortcut

• Repeatedly divide by 2, consider remainder



Read from bottom to top: 100101_B



The Hexadecimal Number System



Name

- "hexa-" (Ancient Greek $\xi\alpha$ -) \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
 - $A13D_{H} \neq 3DA1_{H}$

Computer programmers often use hexadecimal or "hex"

• In C: **0x** prefix (**0xA13D**, etc.)

Why?



Decimal-Hexadecimal Equivalence

				_	
Decimal	Hex	Decimal	Hex		Decimal Hex
0	0	16	10		32 20
1	1	17	11		33 21
2	2	18	12		34 22
3	3	19	13		35 23
4	4	20	14		36 24
5	5	21	15		37 25
6	б	22	16		38 26
7	7	23	17		39 27
8	8	24	18		40 28
9	9	25	19		41 29
10	A	26	1A		42 2A
11	В	27	1B		43 2B
12	C	28	1C		44 2C
13	D	29	1D		45 2D
14	Е	30	1E		46 2E
15	F	31	1F		47 2F

Integer-Hexadecimal Conversion



Hexadecimal to integer: expand using positional notation

$$25_{\rm H} = (2*16^{1}) + (5*16^{0})$$

= 32 + 5
= 37

Integer to hexadecimal: use the shortcut

37 / 16 = 2 R 5 2 / 16 = 0 R 2 Read from bottom to top: 25_H

Binary-Hexadecimal Conversion



Observation: $16^1 = 2^4$

• Every 1 hexadecimal digit corresponds to 4 binary digits

Binary to hexadecimal

1010 0001 0011 1101 _B				
Α	1	3	$D_{_{\mathrm{H}}}$	

Hexadecimal to binary

A 1 3 D_H 1010000100111101_B Digit count in binary number not a multiple of $4 \Rightarrow$ pad with zeros on left

Discard leading zeros from binary number if appropriate

Is it clear why programmers often use hexadecimal?

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

The Octal Number System

Name

• "octo" (Latin) \Rightarrow eight

Characteristics

- Eight symbols
 - 0 1 2 3 4 5 6 7
- Positional
 - 1743₀ ≠ 7314₀

Computer programmers often use octal (so does Mickey!)

• In C: 0 prefix (01743, etc.)





Agenda



Number Systems

Finite representation of unsigned integers

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

Integral Types in Java vs. C

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

- 1. Not guaranteed by C, but on armlab, char = 8 bits, short = 16 bits, int = 32 bits, long = 64 bits
- 2. Not guaranteed by C, but on armlab, char is unsigned

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



Representing Unsigned Integers

Mathematics

Range is 0 to ∞

Computer programming

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

Typical computers today

• n = 32 or 64, so range is 0 to $2^{32} - 1 \text{ or } 2^{64} - 1$ (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

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			

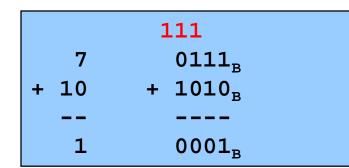
Adding Unsigned Integers



Addition

	1	
3	0011 _B	
+ 10	+ 1010 _B	
13	1101 _B	

Start at right column Proceed leftward Carry 1 when necessary



Results are mod 2⁴



How would you detect overflow programmatically?

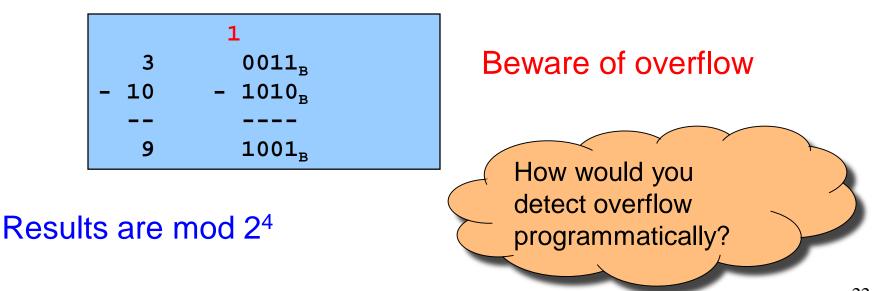




Subtraction

	111	
10	1010 _B	
- 7	- 0111 _B	
3	0011 _B	

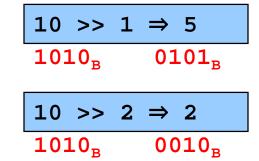
Start at right column Proceed leftward Borrow when necessary



Shifting Unsigned Integers

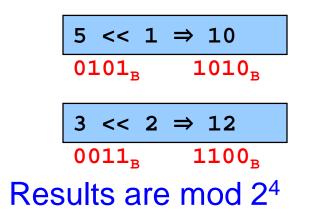


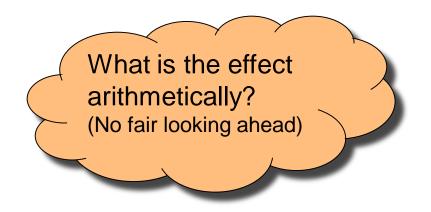
Bitwise right shift (>> in C): fill on left with zeros



What is the effect arithmetically? (No fair looking ahead)

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





Other Operations on Unsigned Ints



Bitwise NOT (~ in C)

• Flip each bit

$$\begin{array}{c} \sim 10 \implies 5 \\ 1010_{\rm B} \qquad 0101_{\rm B} \end{array}$$

Bitwise AND (& in C)

Logical AND corresponding bits

10	1010 _B
& 7	& 0111 _B
2	0010 _B

Useful for setting selected bits to 0

Other Operations on Unsigned Ints



Bitwise OR: (| in C)

Logical OR corresponding bits

	10 1	1010 _B 0001 _B
	11	1011 _B

Useful for setting selected bits to 1

Bitwise exclusive OR (^ in C)

• Logical exclusive OR corresponding bits

10	1010 _B
^ 10	^ 1010 _B
0	0000 _B

x ^ x sets all bits to 0

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);
```

Aside: Using Bitwise Ops for Arith



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

- $x * 2^{y} == x << y$
 - $3*4 = 3*2^2 = 3 << 2 \Rightarrow 12$
- $x / 2^{y} == x >> y$
 - $13/4 = 13/2^2 = 13 >> 2 \Rightarrow 3$

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

•
$$13\%4 = 13\%2^2 = 13\&(2^2-1)$$

13	1101 _B
& 3	& 0011 _B
 1	0001 _B

Fast way to **multiply** by a power of 2

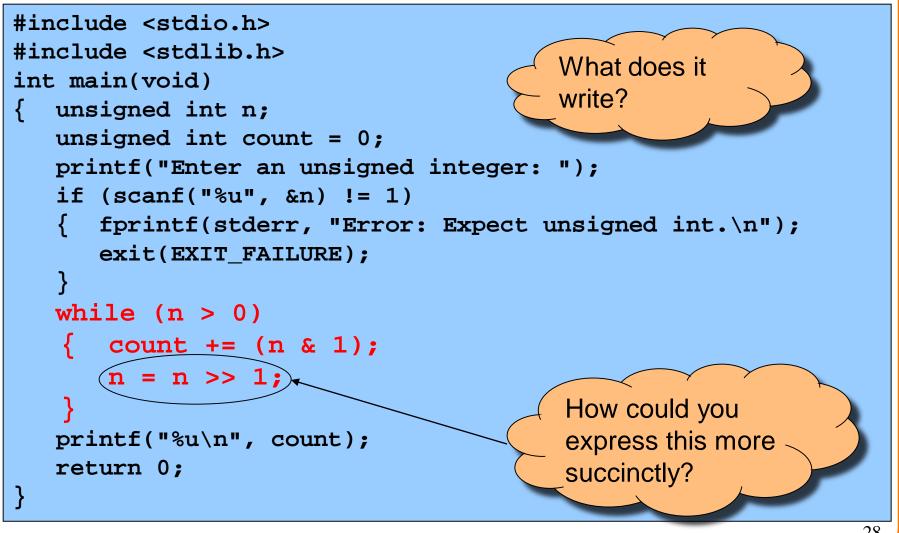
Fast way to **divide <u>unsigned</u>** by power of 2

Fast way to **mod** by a power of 2

Many compilers will do these transformations automatically!

Aside: Example C Program









Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational (floating-point) numbers

Sign-Magnitude



<u>Integer</u>	Rep	
-7	1111	
-6	1110	
-5	1101	
-4	1100	Defir
-3	1011	High-
-2	1010	
-1	1001	0
-0	1000	1
0	0000	Rema
1	0001	
2	0010	0
3	0011	1
4	0100	
5	0101	
6	0110	
7	0111	
		I

Definition High-order bit indicates sign $0 \Rightarrow \text{positive}$ $1 \Rightarrow \text{negative}$ Remaining bits indicate magnitude $0101_B = 101_B = 5$ $1101_B = -101_B = -5$

Sign-Magnitude (cont.)



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

Computing negative neg(x) = flip high order bit of x $neg(0101_B) = 1101_B$ $neg(1101_B) = 0101_B$

Pros and cons

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

Ones' Complement



IntegerRep-71000-61001-51010-41011-31100-21101-11110-011110000010001200103001140100501016011070111	efinition gh-order bit has weight -7 $10_{B} = (1*-7)+(0*4)+(1*2)+(0*1)$ = -5 $10_{B} = (0*-7)+(0*4)+(1*2)+(0*1)$ = 2
--	--

Ones' Complement (cont.)



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

Computing negative neg(x) = -x $neg(0101_B) = 1010_B$ $neg(1010_B) = 0101_B$

Similar pros and cons to sign-magnitude

Two's Complement



IntegerRep-81000-71001-61010-51011-41100-31101-21110-111110000010001200103001140100501016011070111	$\begin{array}{l} \text{hition} \\ \text{-order bit has weight -8} \\ D_{\text{B}} &= (1*-8)+(0*4)+(1*2)+(0*1) \\ &= -6 \\ D_{\text{B}} &= (0*-8)+(0*4)+(1*2)+(0*1) \\ &= 2 \end{array}$
--	--

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 neg(x) = -x + 1 neg(x) = onescomp(x) + 1 $neg(0101_B) = 1010_B + 1 = 1011_B$ $neg(1011_B) = 0100_B + 1 = 0101_B$

Pros and cons

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

Two's Complement (cont.)



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

• Arithmetic is easy!

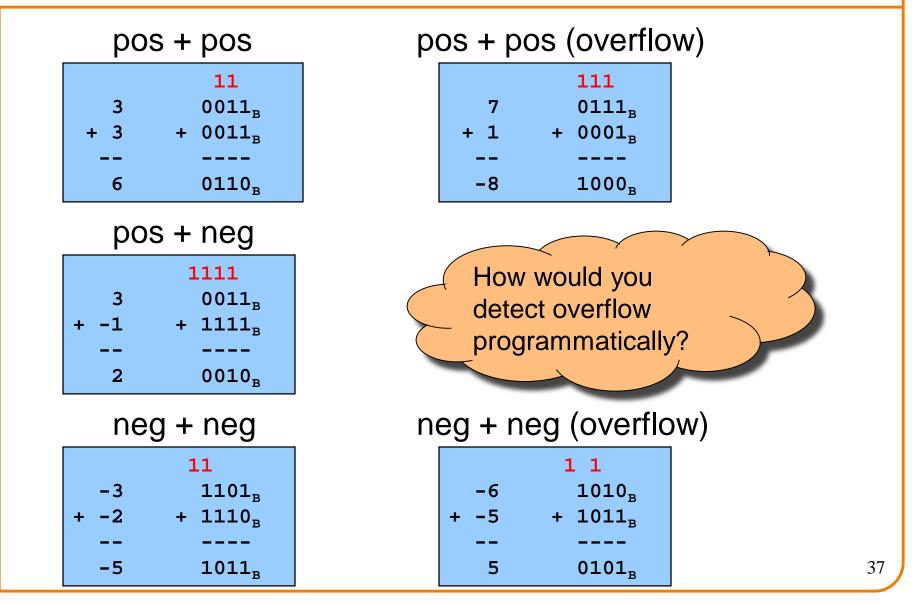
Is it after 1980? OK, then we're surely two's complement



Hereafter, assume two's complement

Adding Signed Integers





Subtracting Signed Integers



Perform subtraction
with borrowsorCompute two's comp
and add11
3
-4
-1
-111
 1111_B 4
-1
-13
 -1111_B

	111
-5	1011
+ -2	+ 1110
-7	1001

Negating Signed Ints: Math



Question: Why does two's comp arithmetic work?

Answer: $[-b] \mod 2^4 = [twoscomp(b)] \mod 2^4$

$$[-b] \mod 2^4$$

$$= [2^4 - b] \mod 2^4$$

$$= [2^4 - 1 - b + 1] \mod 2^4$$

$$= [(2^4 - 1 - b) + 1] \mod 2^4$$

= $[onescomp(b) + 1] \mod 2^4$

= $[twoscomp(b)] \mod 2^4$

See Bryant & O' Hallaron book for much more info





And so:

 $[a - b] \mod 2^4 = [a + twoscomp(b)] \mod 2^4$

$$[a - b] \mod 2^{4}$$

$$= [a + 2^{4} - b] \mod 2^{4}$$

$$= [a + 2^{4} - 1 - b + 1] \mod 2^{4}$$

$$= [a + (2^{4} - 1 - b) + 1] \mod 2^{4}$$

$$= [a + \text{onescomp}(b) + 1] \mod 2^{4}$$

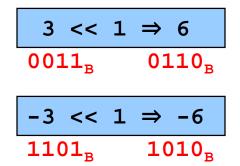
$$= [a + twoscomp(b)] \mod 2^{4}$$

See Bryant & O' Hallaron book for much more info

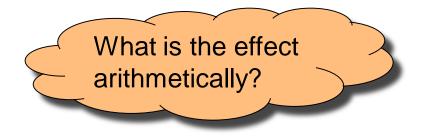
Shifting Signed Integers



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



Results are mod 2⁴



Bitwise right shift: fill on left with ???

Shifting Signed Integers (cont.)



Bitwise arithmetic right shift: fill on left with sign bit



Bitwise logical right shift: fill on left with zeros



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

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



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





Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational (floating-point) numbers

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

Floating Point Numbers



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

This has the form (multiplier) × (base)^(power)

In the computer,

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

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

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

Floating Point Example



Sign (1 bit):

• $1 \Rightarrow$ negative

32-bit representation

Exponent (8 bits):

- $10000011_{B} = 131$
- \cdot 131 127 = 4

Mantissa (23 bits):

- 1 + $(1*2^{-1})+(0*2^{-2})+(1*2^{-3})+(1*2^{-4})+(0*2^{-5})+(1*2^{-6})+(1*2^{-7}) = 1.7109375$

Number:

• $-1.7109375 * 2^4 = -27.375$

When was floating-point invented?



Answer: long before computers!

mantissa

noun

decimal part of a logarithm, 1865, from Latin *mantisa* "a worthless addition, makeweight," perhaps a Gaulish word introduced into Latin via Etruscan (cf. Old Irish *meit*, Welsh *maint* "size").

COM	MON LOGARITHMS

log10x

x 0 I 2 3					100,0000	- 5	6		0		Δ_{993}	Ľ	2	3
	3	4	3	9	7	8	9.	+			Caller .			
50	-6990	6998	7007	7016	7024	7033	7042	7050	7059	7067	9	I	2	1
51	.7076	7084	7093	7101	7110	7118	7126	7135	7143	7152	8	I	2	2
53	.7160		7177		7193	7202	7210		7226	and the second	8	I	2	2
53	.7243	7251	7259	7267	7275	7284	7292	7300	7308	7316	8	I	2	2
54	.7324	7332	7340	7348	7356	7364	7372	7380	7388	7396	8	I	2	2
55	.7404	7412	7419	7427	1 S S S S S S	7443	7451	7450	7466		8	T	2	3

Floating Point Consequences



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

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

- No such number as 1.00000001
- Rule of thumb: "almost 7 digits of precision"

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

• 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	<u>Rational</u>
<u>Approx</u>	<u>Value</u>
.3	3/10
.33	33/100
.333	333/1000
•••	

Binary	<u>Rational</u>
Approx	<u>Value</u>
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
•••	



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

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