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

# Number Systems
## and
# Number Representation

1

## For Your Amusement

**Question**: Why do computer programmers confuse Christmas and Halloween?

**Answer**: Because 25 Dec = 31 Oct

– http://www.electronicsweekly.com

2

## 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 numbers (if time)

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

3

## Agenda

**Number Systems**

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

4

## The Decimal Number System

Name
- "decem" (Latin) ⇒ 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

Why?

5

## 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
  - $1010_B \neq 1100_B$

Most (digital) computers use the binary number system

Terminology
- **Bit**: a binary digit
- **Byte**: (typically) 8 bits

Why?

6

## 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 |
| | | | ... | ... |

7

## Decimal-Binary Conversion

Binary to decimal: expand using positional notation

$$100101_B = (1*2^5) + (0*2^4) + (0*2^3) + (1*2^2) + (0*2^1) + (1*2^0)$$
$$= 32 + 0 + 0 + 4 + 0 + 1$$
$$= 37$$

8

## ~~Decimal~~ Integer-Binary Conversion

Binary to ~~decimal~~ Integer: expand using positional notation

$$100101_B = (1*2^5) + (0*2^4) + (0*2^3) + (1*2^2) + (0*2^1) + (1*2^0)$$
$$= 32 + 0 + 0 + 4 + 0 + 1$$
$$= 37$$

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

9

## Integer-Binary Conversion

Integer to binary: do the reverse
- Determine largest power of 2 ≤ number; write template

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

- Fill in template

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

```
 37  =  (1*2^5)+(0*2^4)+(0*2^3)+(1*2^2)+(0*2^1)+(1*2^0)
-32
  5
 -4
  1                      100101_B
 -1
  0
```

10

## Integer-Binary Conversion

Integer to binary shortcut
- Repeatedly divide by 2, consider remainder

```
37 / 2 = 18 R 1
18 / 2 =  9 R 0
 9 / 2 =  4 R 1
 4 / 2 =  2 R 0
 2 / 2 =  1 R 0
 1 / 2 =  0 R 1
```

Read from bottom to top: $100101_B$

11

## The Hexadecimal Number System

Name
- "hexa" (Greek) ⇒ six
- "decem" (Latin) ⇒ 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 the hexadecimal number system    Why?

12

## 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 | 6 | 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 | B | 27 | 1B | 43 | 2B |
| 12 | C | 28 | 1C | 44 | 2C |
| 13 | D | 29 | 1D | 45 | 2D |
| 14 | E | 30 | 1E | 46 | 2E |
| 15 | F | 31 | 1F | 47 | 2F |
| | | | | ... | ... |

13

## Integer-Hexadecimal Conversion

Hexadecimal to integer: expand using positional notation

```
25ₕ = (2*16¹) + (5*16⁰)
    =  32   +    5
    =  37
```

$$25_H = (2*16^1) + (5*16^0) = 32 + 5 = 37$$

Integer to hexadecimal: use the shortcut

```
37 / 16 = 2 R 5
 2 / 16 = 0 R 2
```

Read from bottom to top: $25_H$

14

## Binary-Hexadecimal Conversion

Observation: $16^1 = 2^4$
- Every 1 hexadecimal digit corresponds to 4 binary digits

Binary to hexadecimal

```
1010 0001 0011 1101ᵦ
 A    1    3    Dₕ
```

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

Hexadecimal to binary

```
 A    1    3    Dₕ
1010 0001 0011 1101ᵦ
```

Discard leading zeros from binary number if appropriate

Is it clear why programmers often use hexadecimal?

15

## The Octal Number System

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

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

Computer programmers often use the octal number system

Why?

(So does Mickey Mouse!)

16

## Agenda

Number Systems

**Finite representation of unsigned integers**

Finite representation of signed integers

Finite representation of rational numbers (if time)

17

## Unsigned Data Types: Java vs. C

Java has type:
- `int`
  - Can represent signed integers

C has type:
- `signed int`
  - Can represent signed integers
- `int`
  - Same as `signed int`
- `unsigned int`
  - Can represent only unsigned integers

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

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**

CourseLab computers
- $n = 64$, so range is 0 to $2^{64} - 1$ (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/subtracting binary numbers

Addition

```
0011
1010
```

Subtraction

```
1010
0111
```

Subtraction

```
0011
1010
```

21

## Adding Unsigned Integers

Addition

```
        1
    3      0011ʙ
+ 10    + 1010ʙ
  --      ----
  13      1101ʙ
```

Start at right column
Proceed leftward
Carry 1 when necessary

```
        1
    7      0111ʙ
+ 10    + 1010ʙ
  --      ----
   1      0001ʙ
```

Beware of overflow

How would you detect overflow programmatically?

Results are mod $2^4$

22

## Subtracting Unsigned Integers

Subtraction

```
       111
   10    1010ʙ
 -  7  - 0111ʙ
   --    ----
    3    0011ʙ
```

Start at right column
Proceed leftward
Borrow when necessary

```
        1
    3    0011ʙ
 - 10  - 1010ʙ
   --    ----
    9    1001ʙ
```

Beware of overflow

How would you detect overflow programmatically?

Results are mod $2^4$

23

## Shifting Unsigned Integers

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

```
10 >> 1 ⇒ 5
1010ʙ   0101ʙ
```

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

```
10 >> 2 ⇒ 2
1010ʙ   0010ʙ
```

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

```
5 << 1 ⇒ 10
0101ʙ   1010ʙ
```

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

```
3 << 2 ⇒ 12
0011ʙ   1100ʙ
```

Results are mod $2^4$

24

## Other Operations on Unsigned Ints

Bitwise NOT (~ in C)
- Flip each bit

```
~10  ⇒ 5
1010B 0101B
```

Bitwise AND (& in C)
- Logical AND corresponding bits

```
  10      1010B
& 7     & 0111B
--        ----
   2      0010B
```

Useful for setting selected bits to 0

25

## Other Operations on Unsigned Ints

Bitwise OR: (| in C)
- Logical OR corresponding bits

```
  10      1010B
| 1     | 0001B
--        ----
  11      1011B
```

Useful for setting selected bits to 1

Bitwise exclusive OR (^ in C)
- Logical exclusive OR corresponding bits

```
  10      1010B
^ 10    ^ 1010B
--        ----
   0      0000B
```

x ^ x sets all bits to 0

26

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

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

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

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

```
  13      1101B
&  3    & 0011B
--        ----
   1      0001B
```

27

## Aside: Example C Program

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

What does it write?

How could this be expressed more succinctly?

28

## Agenda

Number Systems

Finite representation of unsigned integers

**Finite representation of signed integers**

Finite representation of rational numbers (if time)

29

## Signed Magnitude

| 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 |

**Definition**
High-order bit indicates sign
    0 ⇒ positive
    1 ⇒ negative
Remaining bits indicate magnitude
    $1101_B = -101_B = -5$
    $0101_B = 101_B = 5$

30

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

**Pros and cons**
+ easy for people to understand
+ symmetric
- two representations of zero
- can't use the same "add" algorithm for both 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 |
| 0 | 0000 |
| 1 | 0001 |
| 2 | 0010 |
| 3 | 0011 |
| 4 | 0100 |
| 5 | 0101 |
| 6 | 0110 |
| 7 | 0111 |

**Definition**
High-order bit has weight -7
$1010_B = (1*-7)+(0*4)+(1*2)+(0*1) = -5$
$0010_B = (0*-7)+(0*4)+(1*2)+(0*1) = 2$

32

---

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

**Computing negative (alternative)**
neg(x) = $1111_B$ - x
$\text{neg}(0101_B) = 1111_B - 0101_B = 1010_B$
$\text{neg}(1010_B) = 1111_B - 1010_B = 0101_B$

**Pros and cons**
+ symmetric
– two reps of zero
– can't use the same "add" algorithm for both signed and unsigned numbers

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 -8
$1010_B = (1*-8)+(0*4)+(1*2)+(0*1) = -6$
$0010_B = (0*-8)+(0*4)+(1*2)+(0*1) = 2$

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**
neg(x) = ~x + 1
neg(x) = 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
+ one representation of zero
+ same algorithm adds unsigned numbers or signed numbers

35

---

## Two's Complement (cont.)

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

Why?
- Arithmetic is easy
  - Will become clear soon

Hereafter, assume two's complement representation of signed integers

36

## Adding Signed Integers

```
      pos + pos              pos + pos (overflow)
         11                        111
    3     0011ʙ            7     0111ʙ
  + 3   + 0011ʙ          + 1   + 0001ʙ
  --    ----             --    ----
    6     0110ʙ           -8     1000ʙ

      pos + neg
        1111
    3     0011ʙ
  + -1  + 1111ʙ
  --    ----
    2    10010ʙ

      neg + neg              neg + neg (overflow)
        11                      1 1
   -3     1101ʙ           -6     1010ʙ
  + -2  + 1110ʙ          + -5  + 1011ʙ
  --    ----             --    ----
   -5    11011ʙ            5    10101ʙ
```

How would you detect overflow programmatically?

37

## Subtracting Signed Integers

Perform subtraction with borrows    or    Compute two's comp and add

```
        1
       22
   3     0011ʙ                  3     0011ʙ
 - 4   - 0100ʙ               + -4   + 1100ʙ
 --    ----                  --    ----
  -1     1111ʙ                -1     1111ʙ
```

```
  -5     1011ʙ                 -5     1011
 - 2   - 0010ʙ               + -2   + 1110
 --    ----                  --    ----
  -7     1001ʙ                -7    11001
```

38

## Negating Signed Ints: Math

**Question**: Why does two's comp arithmetic work?

**Answer:** `[-b] mod 2⁴ = [twoscomp(b)] mod 2⁴`

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

See Bryant & O'Hallaron book for much more info

39

## Subtracting Signed Ints: Math

**And so**:
`[a - b] mod 2⁴ = [a + twoscomp(b)] mod 2⁴`

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

See Bryant & O'Hallaron book for much more info

40

## Shifting Signed Integers

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

```
3 << 1 ⇒ 6
0011ʙ   0110ʙ
```

```
-3 << 1 ⇒ -6
1101ʙ   -1010ʙ
```

What is the effect arithmetically?

Bitwise **arithmetic** right shift: fill on left **with sign bit**

```
6 >> 1 ⇒ 3
0110ʙ   0011ʙ
```

```
-6 >> 1 ⇒ -3
1010ʙ   1101ʙ
```

Results are mod 2⁴

What is the effect arithmetically?

41

## Shifting Signed Integers (cont.)

Bitwise **logical** right shift: fill on left **with zeros**

```
6 >> 1 ⇒ 3
0110ʙ   0011ʙ
```

```
-6 >> 1 ⇒ 5
1010ʙ   0101ʙ
```

What is the effect arithmetically**???**

In C, right shift (>>) could be logical or arithmetic
 • Not specified by C90 standard
 • Compiler designer decides

**Best to avoid shifting signed integers**

(if you must shift signed integers, make sure you're on a 2's complement machine, and do a bitwise-and after shifting)

(Java does this better, with two operators: >> >>> )

42

## Shifting Signed Integers (cont.)



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

(if you must shift signed integers, <u>make sure you're on a 2's complement machine</u>, and do a bitwise-and after shifting)

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

## Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

**Finite representation of rational numbers (if time)**

45

## Rational Numbers

Mathematics
• A **rational** number is one that can be expressed as the **ratio** of two integers
• Infinite range and precision

Computer science
• Finite range and precision
• Approximate using **floating point** number
  • Binary point "floats" across bits

46

## 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 (unsigned number with bias)
• 23 bits: binary fraction of the form 1.dddddddddddd dd ddd dd ddd

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.dddddddddddd ddd dd ddd dd ddd dd ddd dd ddd dd ddd dd ddd d

47

## Floating Point Example

Sign (1 bit):
• 1 ⇒ negative

1 1 0 0 00 101 1011 011 0000 000 0000 00000

32-bit representation

Exponent representation (8 bits):
• $10000101_B = 133$
• Therefore, exponent $= 133 - 127 = 6$

Fraction (23 bits):     also called "mantissa"
• $1.10110110000000000000000_B$
• $1 + (1*2^{-1}) + (0*2^{-2}) + (1*2^{-3}) + (1*2^{-4}) + (0*2^{-5}) + (1*2^{-6}) + (1*2^{-7}) = 1.7109375$

Number:
• $-1.7109375 * 2^6 = -109.5$

48

## 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").



COMMON LOGARITHMS    log₁₀x

| x | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Δₘ 1 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 |
| 56 | ·7482 | 7490 7497 7505 | 7513 7520 7528 | 7536 7543 7551 | 8 | 1 2 2 |
| 57 | ·7559 | 7566 7574 7582 | 7589 7597 7604 | 7612 7619 7627 | 8 | 1 2 2 |
| 58 | ·7634 | 7642 7649 7657 | 7664 7672 7679 | 7686 7694 7701 | 8 | 1 2 2 |
| 59 | ·7709 | 7716 7723 7731 | 7738 7745 7752 | 7760 7767 7774 | 7 | 1 1 2 |

## Floating Point Warning

Decimal number system can represent only some rational numbers with finite digit count
  • Example: 1/3

| Decimal Approx | Rational Value |
|---|---|
| .3 | 3/10 |
| .33 | 33/100 |
| .333 | 333/1000 |
| ... | |

Binary number system can represent only some rational numbers with finite digit count
  • Example: 1/5

| Binary Approx | Rational Value |
|---|---|
| 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 |
| 0.0011001 1 | 51/256 |
| ... | |

Beware of **roundoff error**
  • Error resulting from inexact representation
  • Can accumulate

50

## Summary

The binary, hexadecimal, and octal number systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers

Essential for proper understanding of
  • C primitive data types
  • Assembly language
  • Machine language

51