



Number Systems and Number Representation

Jennifer Rexford



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
 - $2945 \neq 2495$
 - $2945 = (2*10^3) + (9*10^2) + (4*10^1) + (5*10^0)$

(Most) people use the decimal number system

Why?

5



The Binary Number System

Name

- "binarius" (Latin) => 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
0	0
1	1
2	10
3	11
4	100
5	101
6	110
7	111
8	1000
9	1001
10	1010
11	1011
12	1100
13	1101
14	1110
15	1111
...	...



7

Decimal-Binary Conversion	
Binary to decimal: expand using positional notation	
$ \begin{aligned} 100101_2 &= (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 \end{aligned} $	



8

Decimal-Binary Conversion	
Decimal to binary: do the reverse	
• Determine largest power of 2 \leq number; write template	
$37 = (\underline{?} * 2^5) + (\underline{?} * 2^4) + (\underline{?} * 2^3) + (\underline{?} * 2^2) + (\underline{?} * 2^1) + (\underline{?} * 2^0)$	
• Fill in template	
$ \begin{array}{r} 37 = (1 * 2^5) + (0 * 2^4) + (0 * 2^3) + (\underline{1} * 2^2) + (\underline{0} * 2^1) + (\underline{1} * 2^0) \\ -32 \\ \hline 5 \\ -4 \\ \hline 1 \\ -1 \\ \hline 0 \end{array} \quad 100101_2 $	



9

Decimal-Binary Conversion	
Decimal to binary shortcut	
• Repeatedly divide by 2, consider remainder	
$ \begin{array}{r} 37 / 2 = 18 \text{ R } 1 \\ 18 / 2 = 9 \text{ R } 0 \\ 9 / 2 = 4 \text{ R } 1 \\ 4 / 2 = 2 \text{ R } 0 \\ 2 / 2 = 1 \text{ R } 0 \\ 1 / 2 = 0 \text{ R } 1 \end{array} $	<p style="color: red; margin-left: 20px;">↑ Read from bottom to top: 100101_B</p>



10

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?	



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
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
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

Decimal-Hexadecimal Conversion

Hexadecimal to decimal: expand using positional notation

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

Decimal to hexadecimal: use the shortcut

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

↑ Read from bottom to top: 25_{16}

13

Binary-Hexadecimal Conversion

Observation: $16^1 = 2^4$

- Every 1 hexadecimal digit corresponds to 4 binary digits

Binary to hexadecimal

$$\begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & b \end{array}$$

A 1 3 D₁₆

Hexadecimal to binary

$$\begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & b \end{array}$$

A 1 3 D₁₆

Digit count in binary number not a multiple of 4 => pad with zeros on left

Discard leading zeros from binary number if appropriate

Is it clear why programmers often use hexadecimal?

14

The Octal Number System



Name

- "octo" (Latin) => eight

Characteristics

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

Computer programmers often use the octal number system



15

Decimal-Octal Equivalence



Decimal	Octal
0	0
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	10
9	11
10	12
11	13
12	14
13	15
14	16
15	17

Decimal	Octal
16	20
17	21
18	22
19	23
20	24
21	25
22	26
23	27
24	30
25	31
26	32
27	33
28	34
29	35
30	36
31	37

Decimal	Octal
32	40
33	41
34	42
35	43
36	44
37	45
38	46
39	47
40	50
41	51
42	52
43	53
44	54
45	55
46	56
47	57

16

Decimal-Octal Conversion



Octal to decimal: expand using positional notation

$$\begin{aligned} 37_8 &= (3 * 8^1) + (7 * 8^0) \\ &= 24 + 7 \\ &= 31 \end{aligned}$$

Decimal to octal: use the shortcut

$$\begin{array}{r} 31 / 8 = 3 \text{ R } 7 \\ 3 / 8 = 0 \text{ R } 3 \end{array}$$

↑ Read from bottom to top: 37_8

17

Binary-Octal Conversion



Observation: $8^1 = 2^3$

- Every 1 octal digit corresponds to 3 binary digits

Binary to octal

$$\begin{array}{cccc} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & b \end{array}$$

001 010 000 1001 111 01b
1 2 0 4 7 5₈

Digit count in binary number not a multiple of 3 => pad with zeros on left

Octal to binary

$$\begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & b \end{array}$$

1 2 0 4 7 5₈
001 010 000 1001 111 01b

Discard leading zeros from binary number if appropriate

Is it clear why programmers sometimes use octal?

18

Agenda



Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

19

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

20

Representing Unsigned Integers



Mathematics

- Range is 0 to ∞

Computer programming

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

FC010 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

21

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

22

Adding Unsigned Integers



Addition

$$\begin{array}{r} & 1 \\ & 0011_2 \\ \text{3} & + 1010_2 \\ + 10 & \hline \text{---} \\ 13 & 1101_2 \end{array}$$

Start at right column
Proceed leftward
Carry 1 when necessary

$$\begin{array}{r} & 11 \\ & 0111_2 \\ \text{7} & + 1010_2 \\ + 10 & \hline \text{---} \\ 1 & 0001_2 \end{array}$$

Beware of overflow

Results are mod 2^4

How would you detect overflow programmatically?

23

Subtracting Unsigned Integers



Subtraction

$$\begin{array}{r} & 12 \\ & 0202 \\ \text{10} & - 0111_2 \\ - 7 & \hline \text{---} \\ 3 & 0011_2 \end{array}$$

Start at right column
Proceed leftward
Borrow 2 when necessary

$$\begin{array}{r} & 2 \\ & 0011_2 \\ \text{3} & - 1010_2 \\ - 10 & \hline \text{---} \\ 9 & 1001_2 \end{array}$$

Beware of overflow

Results are mod 2^4

How would you detect overflow programmatically?

24

Shifting Unsigned Integers



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

$10 >> 1 \Rightarrow 5$
$1010_8 \quad 0101_8$
$10 >> 2 \Rightarrow 2$
$1010_8 \quad 0010_8$

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

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

$5 << 1 \Rightarrow 10$
$0101_8 \quad 1010_8$
$3 << 2 \Rightarrow 12$
$0011_8 \quad 1100_8$

Results are mod 2^4

25

Other Operations on Unsigned Ints



Bitwise NOT (`~` in C)

- Flip each bit

$\sim 10 \Rightarrow 5$
$1010_8 \quad 0101_8$

Bitwise AND (`&` in C)

- Logical AND corresponding bits

10	1010_8
$\& 7$	$\& 0111_8$
$--$	$-----$
2	0010_8

Useful for setting selected bits to 0

26

Other Operations on Unsigned Ints



Bitwise OR: (`|` in C)

- Logical OR corresponding bits

10	1010_8
$ 1$	$ 0001_8$
$--$	$-----$
11	1011_8

Useful for setting selected bits to 1

Bitwise exclusive OR (`^` in C)

- Logical exclusive OR corresponding bits

10	1010_8
$\wedge 10$	$\wedge 1010_8$
$--$	$-----$
0	0000_8

$x \wedge x$ sets all bits to 0

27

Aside: Using Bitwise Ops for Arith



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

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

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

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

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

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

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

13	1101_8
$\& 3$	$\& 0011_8$
$--$	$-----$
1	0001_8

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

28

Aside: Example C Program



```
#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?

29

Agenda



Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

30

Signed 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

$$1101_B = -101_B = -5$$

$$0101_B = 101_B = 5$$

31

Signed 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 for people to understand

+ symmetric

- two reps of zero

32

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	0011
4	0100
5	0101
6	0110
7	0111

Definition

High-order bit has weight -7

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

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

33

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	0011
4	0100
5	0101
6	0110
7	0111

Computing negative

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

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

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

Computing negative (alternative)

$\text{neg}(x) = 1111_B - x$

$$\text{neg}(0101_B) = 1111_B - 0101_B = 1010_B$$

$$\text{neg}(1010_B) = 1111_B - 1010_B = 0101_B$$

Pros and cons

+ symmetric

- two reps of zero

34

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

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

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

35

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

+ one rep of zero

36

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

37

Adding Signed Integers



pos + pos

3	0011 _b	11
+ 3	+ 0011 _b	
--	---	
6	0110 _b	

pos + neg

3	0011 _b	1111
+ -1	+ 1111 _b	
--	---	
2	10010 _b	

neg + neg

-3	1101 _b	11
+ -2	+ 1110 _b	
--	---	
-5	11011 _b	

pos + pos (overflow)

7	0111 _b	111
+ 1	+ 0001 _b	
--	---	
-8	1000 _b	

neg + neg (overflow)

-6	1010 _b	1 1
+ -5	+ 1011 _b	
--	---	
5	10101 _b	

How would you detect overflow programmatically?

38

Subtracting Signed Integers



Perform subtraction with borrows

or

Compute two's comp and add

3	0011 _b	1
- 4	- 0100 _b	22
--	---	
-1	1111 _b	



3	0011 _b	3
+ -4	+ 1100 _b	
--	---	
-1	1111 _b	

-5	1011 _b	-
- 2	- 0010 _b	22
--	---	
-7	1001 _b	



-5	1011 _b	111
+ -2	+ 1110 _b	
--	---	
-7	11001 _b	

39

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 & O'Hallaron book for much more info

40

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 & O'Hallaron book for much more info

41

Shifting Signed Integers



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

3	0011 _b	00110
-3	-1101 _b	-1010

What is the effect arithmetically?

Bitwise arithmetic right shift: fill on left with sign bit

6	0110 _b	0110
-6	-1010 _b	-1101

What is the effect arithmetically?

Results are mod 2⁴

42

Shifting Signed Integers (cont.)

Bitwise logical right shift: fill on left with zeros

6 >> 1 => 3
0110 _b 0011 _b

-6 >> 1 => 5
1010 _b 0101 _b

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

43

The logo of the University of Washington, featuring a crest with a shield containing a tree, surrounded by the words "UNIVERSITY OF WASHINGTON".

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

The logo of the University of California, Berkeley, featuring a shield with a star, an oak tree, and the text "BERKELEY" and "UNIVERSITY OF CALIFORNIA".

Agenda

Number Systems

Finite representation of unsigned integers

Finite representation of signed integers

Finite representation of rational numbers (if time)

The logo of the University of Colorado Boulder, featuring a shield with a mountain peak, a river, and a sun, with the text "UNIVERSITY OF COLORADO BOULDER" around it.

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

Floating Point Warning



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

- Example: 1/3

Decimal	Rational
Approx	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	Rational
Approx	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
0.00110011	51/256
...	

Beware of roundoff error

- Error resulting from inexact representation
- Can accumulate

49

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

50