

# Bits, Bytes, and Integers

Kai Shen

1

## Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating

2

## Binary Representations

3


## Encoding Byte Values

- Byte = 8 bits
  - Binary  $00000000_2$  to  $11111111_2$
  - Decimal:  $0_{10}$  to  $255_{10}$
  - Hexadecimal  $00_{16}$  to  $FF_{16}$ 
    - Base 16 number representation
    - Use characters '0' to '9' and 'A' to 'F'
    - Write  $FA1D37B_{16}$  in C as
      - $0xFA1D37B$
      - $0xfa1d37b$

	Hex	Decimal	Binary
0	0	0000	
1	1	0001	
2	2	0010	
3	3	0011	
4	4	0100	
5	5	0101	
6	6	0110	
7	7	0111	
8	8	1000	
9	9	1001	
A	10	1010	
B	11	1011	
C	12	1100	
D	13	1101	
E	14	1110	
F	15	1111	

4

### Byte-Oriented Memory Organization



- Programs Refer to Virtual Addresses
  - Conceptually very large array of bytes
  - Actually implemented with hierarchy of different memory types
  - System provides address space private to particular "process"
    - Program being executed
    - Program can clobber its own data, but not that of others
- Compiler + Run-Time System Control Allocation
  - Where different program objects should be stored
  - All allocation within single virtual address space

5

### Machine Words

- Machine Has "Word Size"
  - Nominal size of integer-valued data
    - Including addresses
  - Most current machines use 32 bits (4 bytes) words
    - Limits addresses to 4GB
    - Becoming too small for memory-intensive applications
  - High-end systems use 64 bits (8 bytes) words
    - Potential address space  $\approx 1.8 \times 10^{19}$  bytes
    - x86-64 machines support 48-bit addresses: 256 Terabytes
  - Machines support multiple data formats
    - Fractions or multiples of word size
    - Always one or multiples of bytes

6

### Word-Oriented Memory Organization

- Addresses Specify Byte Locations
  - Address of first byte in word
  - Addresses of successive words differ by 4 (32-bit) or 8 (64-bit)

32-bit Words	64-bit Words	Bytes	Addr.
			0000
Addr = 0000			0001
			0002
	Addr = 0000		0003
			0004
Addr = 0004			0005
			0006
			0007
			0008
			0009
Addr = 0008			0010
	Addr = 0008		0011
			0012
Addr = 0012			0013
			0014
			0015

7

### Data Representations

C Data Type	Typical 32-bit	Intel IA32	x86-64
char	1	1	1
short	2	2	2
int	4	4	4
long	4	4	8
long long	8	8	8
float	4	4	4
double	8	8	8
long double	8	10/12	10/16
pointer	4	4	8

8

### Byte Ordering

- How should bytes within a multi-byte word be ordered in memory?
- Big Endian (Sun, PPC Mac, Internet)
  - Most significant byte has lowest address
- Little Endian (Intel x86)
  - Least significant byte has lowest address
- Example
  - Variable x has 4-byte representation 0x01234567
  - Address given by &x is 0x100

**Big Endian**

0x100	0x101	0x102	0x103
01	23	45	67

**Little Endian**

0x100	0x101	0x102	0x103
67	45	23	01

9

### Reading Byte-Reversed Listings

- Disassembly
  - Text representation of binary machine code
  - Generated by program that reads the machine code
- Example Fragment
 

Address	Instruction Code	Assembly Rendition
8048365:	5b	pop %ebx
8048366:	81 c3 ab 12 00 00	add \$0x12ab,%ebx
804836c:	83 bb 28 00 00 00	cmpl \$0x0,0x28(%ebx)
- Deciphering Numbers
  - Value: 0x12ab
  - Pad to 32 bits: 0x00012ab
  - Split into bytes: 00 00 12 ab
  - Reverse: ab 12 00 00

10

### Representing Integers

Decimal: 15213  
Binary: 0011 1011 0110 1101  
Hex: 3 B 6 D

int A = 15213;

IA32, x86-64	Sun
6D	00
3B	00
00	3B
00	6D

long int C = 15213;

IA32	x86-64	Sun
6D	6D	00
3B	3B	00
00	00	3B
00	00	6D
00	00	
00	00	
00	00	

int B = -15213;

IA32, x86-64	Sun
93	FF
C4	FF
FF	C4
FF	93

Two's complement representation (Covered later)

11

### Representing Strings

char S[6] = "18243";

- Strings in C
  - Represented by array of characters
  - Each character encoded in ASCII
    - Standard 7-bit encoding of characters
    - Character "0" has code 0x30
      - Digit i has code 0x30+i
  - String should be null-terminated
    - Final character = 0
- Compatibility
  - Byte ordering not an issue

Linux/Alpha	Sun
31	31
38	38
32	32
34	34
33	33
00	00

12

## Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating

13

## Boolean Algebra

- Algebraic representation of logic
  - Encode "True" as 1 and "False" as 0

And

■  $A \& B = 1$  when both  $A=1$  and  $B=1$

&	0	1
0	0	0
1	0	0

Or

■  $A | B = 1$  when either  $A=1$  or  $B=1$

	0	1
0	0	1
1	1	1

Not

■  $\sim A = 1$  when  $A=0$

~	0	1
0	1	0
1	0	0

Exclusive-Or (Xor)

■  $A \wedge B = 1$  when either  $A=1$  or  $B=1$ , but not both

^	0	1
0	0	1
1	1	0

14

## Boolean Algebra on Bit Vectors

- Operate on Bit Vectors
  - Operations applied bitwise

01101001	01101001	01101001	01010101
& 01010101	01010101	^ 01010101	~ 01010101
01000001	01111101	00111100	10101010

15

## Representing & Manipulating Sets

- Representation
  - Width  $w$  bit vector represents subsets of  $\{0, \dots, w-1\}$
  - $a_j = 1$  if  $j \in A$ 
    - 01101001      $\{0, 3, 5, 6\}$
    - 76543210
    - 01010101      $\{0, 2, 4, 6\}$
    - 76543210
- Operations
 

■ & Intersection	01000001	$\{0, 6\}$
■   Union	01111101	$\{0, 2, 3, 4, 5, 6\}$
■ ^ Symmetric difference	00111100	$\{2, 3, 4, 5\}$
■ ~ Complement	10101010	$\{1, 3, 5, 7\}$

16

## Bit-Level Operations in C

- Operations `&`, `|`, `~`, `^` Available in C
  - Apply to any “integral” data type
    - long, int, short, char, unsigned
  - View arguments as bit vectors
  - Arguments applied bit-wise
- Examples
  - `~0x41 = 0xBE`
    - `~010000012 = 101111102`
  - `~0x00 = 0xFF`
    - `~000000002 = 111111112`
  - `0x69 & 0x55 = 0x41`
    - `011010012 & 010101012 = 010000012`
  - `0x69 | 0x55 = 0x7D`
    - `011010012 | 010101012 = 011111012`

17

## Contrast: Logic Operations in C

- Contrast to Logical Operators
  - `&&`, `||`, `!`
    - View 0 as “False”
    - Anything nonzero as “True”
    - Always return 0 or 1
    - **Early termination**
- Examples (char data type)
  - `!0x41 = 0x00`
  - `!0x00 = 0x01`
  - `!!0x41 = 0x01`
  - `0x69 && 0x55 = 0x01`
  - `0x69 || 0x55 = 0x01`
  - `p && *p` (avoids null pointer access)

18

## Shift Operations

- Left Shift: `x << y`
  - Shift bit-vector `x` left `y` positions
    - Throw away extra bits on left
    - Fill with 0's on right
- Right Shift: `x >> y`
  - Shift bit-vector `x` right `y` positions
    - Throw away extra bits on right
  - Logical shift
    - Fill with 0's on left
  - Arithmetic shift
    - Replicate most significant bit on right

Argument <code>x</code>	01100010
<code>&lt;&lt; 3</code>	00010000
Log. <code>&gt;&gt; 2</code>	00011000
Arith. <code>&gt;&gt; 2</code>	00011000

Argument <code>x</code>	10100010
<code>&lt;&lt; 3</code>	00010000
Log. <code>&gt;&gt; 2</code>	00101000
Arith. <code>&gt;&gt; 2</code>	11101000

19

## Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating

20

### Encoding Integers

Unsigned

$$B2U(X) = \sum_{i=0}^{w-1} x_i \cdot 2^i$$

Two's Complement

$$B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$$

```
short int x = 15213;
short int y = -15213;
```

Sign Bit

- C short 2 bytes long

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
y	-15213	C4 93	11000100 10010011

- Sign Bit
  - For 2's complement, most significant bit indicates sign
    - 0 for nonnegative, 1 for negative
- What is the bit representation of -1?

21

### Encoding Example (Cont.)

```
x = 15213: 00111011 01101101
y = -15213: 11000100 10010011
```

Weight	15213	-15213
1	1	1
2	0	0
4	1	4
8	1	8
16	0	0
32	1	32
64	1	64
128	0	0
256	1	256
512	1	512
1024	0	0
2048	1	2048
4096	1	4096
8192	1	8192
16384	0	0
-32768	0	0
<b>Sum</b>	<b>15213</b>	<b>-15213</b>

22

### Numeric Ranges

Unsigned Values

- $UMin = 0$   
000...0
- $UMax = 2^w - 1$   
111...1

Two's Complement Values

- $TMin = -2^{w-1}$   
100...0
- $TMax = 2^{w-1} - 1$   
011...1

Values for  $W = 16$

	Decimal	Hex	Binary
UMax	65535	FF FF	11111111 11111111
TMax	32767	7F FF	01111111 11111111
TMin	-32768	80 00	10000000 00000000
-1	-1	FF FF	11111111 11111111
0	0	00 00	00000000 00000000

23

### Values for Different Word Sizes

	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

Observations

- $|TMin| = TMax + 1$ 
  - Asymmetric range
- $UMax = 2 * TMax + 1$

C Programming

- #include <limits.h>
- Declares constants, e.g.,
  - ULONG\_MAX
  - LONG\_MAX
  - LONG\_MIN
- Values platform specific

24

### Unsigned & Signed Numeric Values

X	B2U(X)	B2T(X)
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	-8
1001	9	-7
1010	10	-6
1011	11	-5
1100	12	-4
1101	13	-3
1110	14	-2
1111	15	-1

- Equivalence
  - Same encodings for nonnegative values
- Uniqueness
  - Every bit pattern represents unique integer value
  - Each representable integer has unique bit encoding

⇒ Can Invert Mappings

- $U2B(x) = B2U^{-1}(x)$ 
  - Bit pattern for unsigned integer
- $T2B(x) = B2T^{-1}(x)$ 
  - Bit pattern for two's comp integer

25

### Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating

26

### Mapping Between Signed & Unsigned

Two's Complement  $x$  → [T2B] →  $X$  → [B2U] → Unsigned  $ux$   
 Maintain Same Bit Pattern

Unsigned  $ux$  → [U2B] →  $X$  → [B2T] → Two's Complement  $x$   
 Maintain Same Bit Pattern

- Mappings between unsigned and two's complement numbers: **keep bit representations and reinterpret**

27

### Mapping Signed ↔ Unsigned

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

↔ =

↔ +/- 16

28

### Relation between Signed & Unsigned

Two's Complement  $x$  → **T2U** → Unsigned  $ux$

Maintain Same Bit Pattern

Large negative weight becomes large positive weight

$$ux = \begin{cases} x & x \geq 0 \\ x + 2^w & x < 0 \end{cases}$$

29

### Conversion Visualized

- 2's Comp. → Unsigned
- Ordering Inversion
- Negative → Big Positive

2's Complement Range

Unsigned Range

30

### Signed vs. Unsigned in C

- Constants
  - By default are considered to be signed integers
  - Unsigned if have "U" as suffix  
0U, 4294967295U
- Casting
  - Explicit casting between signed & unsigned same as U2T and T2U  

```
int tx, ty;
unsigned ux, uy;
tx = (int) ux;
uy = (unsigned) ty;
```
  - Implicit casting also occurs via assignments and procedure calls  

```
tx = ux;
uy = ty;
```

31

### Casting Surprises

- Expression Evaluation
  - If there is a mix of unsigned and signed in single expression, **signed values implicitly cast to unsigned**
  - Examples for  $W = 32$ :  $TMIN = -2,147,483,648$ ,  $TMAX = 2,147,483,647$

Constant <sub>1</sub>	Constant <sub>2</sub>	Relation	Evaluation
0	0U	==	unsigned
-1	0	<	signed
-1	0U	>	unsigned
2147483647	-2147483647-1	>	signed
2147483647U	-2147483647-1	<	unsigned
-1	-2	>	signed
(unsigned)-1	-2	>	unsigned
2147483647	2147483648U	<	signed
2147483647	(int) 2147483648U	>	signed

32

## Security Vulnerability

```

/* Kernel memory region holding user-accessible data */
#define KSIZE 1024
char kbuf[KSIZE];

/* Copy at most maxlen bytes from kernel region to user buffer */
int copy_from_kernel(void *user_dest, int maxlen) {
    /* Byte count len is minimum of buffer size and maxlen */
    int len = KSIZE < maxlen ? KSIZE : maxlen;
    memcpy(user_dest, kbuf, len);
    return len;
}

```

```

#define MSIZE 528

void getstuff() {
    char mybuf[MSIZE];
    copy_from_kernel(mybuf, MSIZE);
    printf("%s\n", mybuf);
}

```

- Real example: similar to code found in FreeBSD's implementation of getpeername

1/24/2012

CSC252 - Spring 2012

33

## Malicious Usage

```

/* Kernel memory region holding user-accessible data */
#define KSIZE 1024
char kbuf[KSIZE];

/* Copy at most maxlen bytes from kernel region to user buffer */
int copy_from_kernel(void *user_dest, int maxlen) {
    /* Byte count len is minimum of buffer size and maxlen */
    int len = KSIZE < maxlen ? KSIZE : maxlen;
    memcpy(user_dest, kbuf, len);
    return len;
}

```

```

#define MSIZE 528

void getstuff() {
    char mybuf[MSIZE];
    copy_from_kernel(mybuf, -MSIZE);
    . . .
}

```

- MSIZE, when interpreted as unsigned int by memcpy, becomes a very large integer

1/24/2012

CSC252 - Spring 2012

34

## Casting Signed ↔ Unsigned: Basic Rules

- Bit pattern is maintained
- But reinterpreted
- Can have unexpected effects: adding or subtracting  $2^w$
- Expression containing signed and unsigned int
  - int is cast to unsigned!!

35

## Bits, Bytes, and Integers

- Representing information as bits
- Bit-level manipulations
- Integers
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating

36

### Sign 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:
  - $X' = X_{w-1}, \dots, X_{w-1}, X_{w-1}, X_{w-2}, \dots, X_0$

37

### Sign Extension Example

```

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

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
ix	15213	00 00 3B 6D	00000000 00000000 00111011 01101101
y	-15213	C4 93	11000100 10010011
iy	-15213	FF FF C4 93	11111111 11111111 11000100 10010011

- Converting from smaller to larger integer data type
- C automatically performs sign extension

38

### Expanding, Truncating: Basic Rules

- Expanding (e.g., short int to int)
  - Unsigned: zeros added
  - Signed: sign extension
  - Both yield expected result
- Truncating (e.g., unsigned to unsigned short)
  - Unsigned/signed: bits are truncated
  - Result reinterpreted
  - Unsigned: mod operation
  - Signed: similar to mod
  - For small numbers yields expected behavior

39

### Disclaimer

These slides were adapted from the CMU course slides provided along with the textbook of "Computer Systems: A programmer's Perspective" by Bryant and O'Hallaron. The slides are intended for the sole purpose of teaching the computer organization course at the University of Rochester.

40