

Data Representation and Manipulation.

①
- humans - decimal (base 10) number system

- computers - binary (base 2) " " " "

∴ signal - high or low, on or off.

3 representations of numbers

1. unsigned encodings (for numbers ≥ 0)

2. 2's complement " (for signed numbers)

3. floating-point " (for real numbers).

INFORMATION STORAGE

* Hexadecimal Notation

1 byte = 8 bits

$\frac{00000000}{\downarrow \downarrow} 2 \rightarrow 00$ to $\frac{11111111}{\downarrow \downarrow} 2 \rightarrow FF_{16}$

0x or 0X

Symbols: 0, 1, 2, ..., 9, A, B, C, D, E, F.

eg. 0X1A0F = 0001 1010 0000 1111₂

* Words

word size (w) = 32 bits (in most machines).

if w = 32 bits, a program's virtual address space can access addresses in the range of 0 to $2^{32}-1$ (ie) 2^{32} bytes.

w = 64 bits, common these days.

2^w bytes.

* Data sizes

(2)

C type	32-bit	64-bit
char	1	1
short int	2	2
int	4	4
long int	4	8
char *	4	8
float	4	4
double	8	8

short int - at least 16 bits

int - " 16 "

long int = " 32 "

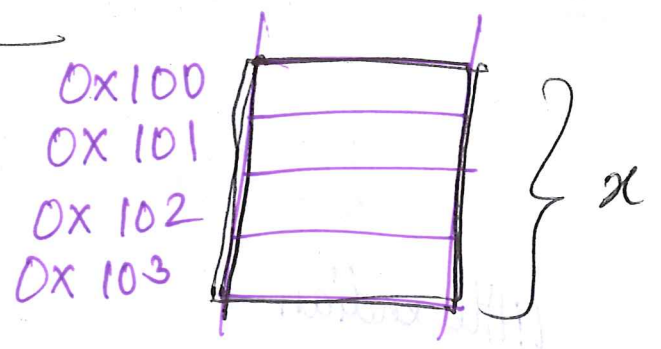
long long int = " 64 "

	Win 16 API	Unix & Unix-like
int	2 bytes	4 bytes.

* Addressing and Byte Ordering

```
int x;
```

```
x = 0x 12345678
```



Endian-ness

Little Endian

0x100	78
101	56
102	34
103	12

Least Significant Byte (LSB) comes first. "most" Intel machines

Big Endian

0x100	12
101	34
102	56
103	78

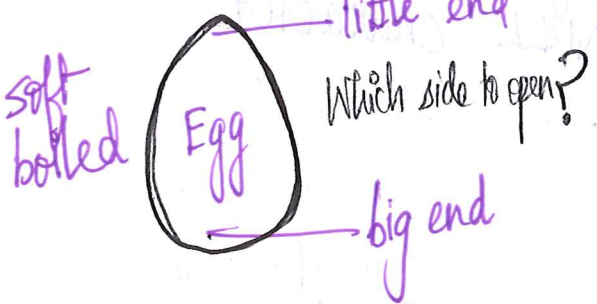
MSB comes first. "most" IBM, Sun Microsystems machines.

⊗ Show code to determine endian-ness of a machine.

Bi-endian - machines that can be configured to work as both little or big endian.

Origin of term "endian-ness"

Gulliver's Travels by Jonathan Swift.



Lilliput and Blefuscu. On holy wars and a plea for peace. (1981) - Danny Cohen. (Networking Protocol).

To see ASCII code: man ascii
 Unicode standard for text:

- ASCII - only English
- Unicode - Albanian to Xamtanga.

lang. spoken by Xamir people in Ethiopia.

↓
 Universal Character Set - 100,000 characters - 32-bit representation of chars.

- UTF-8 - Universal Coded Character Set Transformation Format - 8 bit.
- 1,112,064 valid code points.
- Java uses unicode.

BOOLEAN ALGEBRA

- George Boole around 1850.

TRUE - 1
 FALSE - 0

op.	Logical op.	boolean op.
NOT	$\neg P$	$\sim P$
AND	$P \wedge Q$	$P \& Q$
OR	$P \vee Q$	$P Q$
EXOR	$P \oplus Q$	$P \wedge Q$

(6)

Bit Vectors

strings of 0s and 1s with some fixed length w .

$$a = 10101$$

$$b = \underline{01100}$$

$$a \& b = \underline{00100}$$

Bit-level Operations in C

$$\sim 0x00 \quad \sim [0000 \ 0000] = [1111 \ 1111] \quad (0x) \quad 0xFF$$

$$0x01 \mid 0x02 = \begin{array}{r} 0000 \ 0001 \\ 0000 \ 0010 \\ \hline 0000 \ 0011 \end{array} = 0x03$$

Masking Operations

$$x = 0x \ 89 \ AB \ CD \ EF \quad x \ \& \ \underline{0xFF} = 0x \ 00 \ 00 \ 00 \ EF$$

↓
Masking of the LSB.

* Logical Operations in C (17)

OR - ||

AND - &&

NOT - !

TRUE - any non-zero value

FALSE - 0.

eg. !0x41 = 0x00

!0x00 = 0x01

0x01 && 0x02 = 0x01

0x01 || 0x03 = 0x01

* short circuit in logical operations.

a && 5/a \Rightarrow never cause a division by zero.

p && *p++ \Rightarrow " " " null ptr dereference

* Shift operations in C

x = 01100011

x << 4 = 0011 0000

x >> 4 = 0000 0110
(logical)

x >> 4 = ~~0000~~ 0000 0110
(arithmetic)

y = 1001 0101

y << 4 = 0101 0000

y >> 4 = 0000 1001
(logical)

y >> 4 = 1111 1001
arithmetic

In C, no precise definition for right shifts.

In Java, x >> k means arithmetic and x >>> k means logical.

Integer Representations (8)

$$B2U_w(\vec{x}) = \sum_{i=0}^{w-1} x_i 2^i$$

I. Unsigned encoding.

$$\frac{0}{2^3} \frac{1}{2^2} \frac{0}{2^1} \frac{1}{2^0} = 1 \times 2^2 + 1 \times 2^0 = 5.$$

w=4

II. 2's - Complement Encodings

$$B2T_w(\vec{x}) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$$

w=4

$$0000 = 0$$
$$0001 = 1$$

w=2

$$00 = 0$$
$$01 = 1$$
$$10 = -2$$
$$11 = -1$$

$$B2T_4([1011]) = -1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$
$$= -8 + 2 + 1 = \underline{\underline{-5}}$$

w=3

$$000 = 0$$
$$001 = 1$$
$$010 = 2$$
$$011 = 3$$
$$100 = -4$$
$$101 = -3$$
$$110 = -2$$
$$111 = -1$$

Note: 1's complement ⁽⁹⁾ encodings

ASIDE

$$000 = +0$$

$$001 = +1$$

$$010 = +2$$

$$011 = +3$$

$$100 = -1$$

$$101 = -2$$

$$110 = -3$$

$$111 = -0$$

$$\binom{3-1}{-2 \ -1} \text{ to}$$

$$\binom{3-1}{2 \ -1}$$

Important Numbers

Word size $w = \underline{32 \text{ bits}}$

unsigned \leftarrow $UMax_{32} = 0x \text{ FF FF FF FF} = 4,294,967,295$

$$TMax_{32} = 0x \text{ 7F FF FF FF} = 2,147,483,647$$

Two's complement \leftarrow

$$TMin_{32} = 0x \text{ 80 00 00 00} = -2,147,483,648$$

$$0 = 0x \text{ 00 00 00 00} = 0$$

$$-1 = 0x \text{ FF FF FF FF} = -1$$

$$|TMin| = |TMax| + 1$$

\Rightarrow There is no positive counterpart for TMin.

$$UMax = 2TMax + 1$$

$\langle \text{limits.h} \rangle$ - defines constants with the limits of fundamental integral types for the specific system and compiler implementation used. (10)

eg. $\text{INT_MAX} = +2,147,483,647$
 $\text{INT_MIN} = -2,147,483,648$
 $\text{UINT_MAX} = 4,294,967,295$

} on 32-bit
unix-like machines

ISO C99 standard - `stdint.h`

intN_t and uintN_t

eg. int32_t

uint32_t

MACROS: INTN_MIN , INTN_MAX , and UINTN_MAX

eg. INT32_MIN

* Signed-Magnitude form

$$B2S_w(\vec{x}) = (-1)^{x_{w-1}} \cdot \left(\sum_{i=0}^{w-2} x_i 2^i \right)$$

00 = +0

01 = +1

10 = -0

11 = -1

Two's complement of $x = 2^w - x$ (a single two).
(or)
 $-x$ in two's complement

Ones' complement of $x = [1111 \dots 1] - x$ (multiple ones).
(or)
 $-x$ in ones' complement

* Conversions between Signed and Unsigned

unsigned $u = 4,294,967,295u;$ /* U Max 32 */

int $tu = (int) u;$

printf("u = %u , tu = %d \n", u, tu);

$u = 4294967295$, $tu = -1$

u and tu in hex = $0x\text{FF}\text{FF}\text{FF}\text{FF}$

- when signed \longleftrightarrow unsigned (of same word size)
 \Rightarrow numeric values might change
but
the bit patterns do not.

Signed vs unsigned in ⁽¹²⁾C

- most numbers are signed by default.
eg. 12345 or 0x1A2B

- unsigned: 12345U or 0x1A2Bu.

if signed and unsigned appear in an operation,
signed \rightarrow unsigned.

problems:

$-1 < 0U \Rightarrow 4,294,967,295U < 0U$

Sign-extension
&
Zero-extension

short sx = -12345; **FALSE** X

unsigned short usx = sx; /* 53191 */

int x = sx; /* -12345 */

unsigned ux = usx; /* 53191 */

sx = -12345 : CF C7

usx = 53191 : CF C7

x = -12345 : **FF FF** CF C7

ux = 53191 : **00 00** CF C7

\rightarrow sign extension.

\rightarrow zero extension.

sign extension

w=3 to w=4

$$[101]$$

$$-4 + 1 = -3$$

$$[1101]$$

$$-8 + 4 + 1 = -3$$

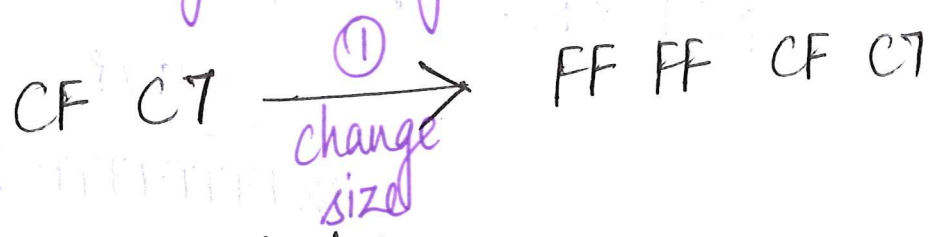
$$[1111] = [111] = -1$$

short to unsigned int ?

short SX = -12345;

/* CF CF */

unsigned uy = SX;



② change from signed to unsigned.

(ie.) interpreted as 429495495 ✓

if 1. signed \rightarrow unsigned } \Rightarrow 53,191. X
2. change size

* Truncating Numbers

int x = 53191;

short sx = (short) x; /* -12345 */

int y = sx; /* -12345 */
(\because sign extension).

Advice on Signed vs. Unsigned ⁽¹⁴⁾

```
size_t strlen(const char *s);
```

↓
unsigned int

```
int strcmp(const char *s, const char *t)
```

```
{ return strlen(s) - strlen(t) > 0;
```

```
}
```

When

$s > t$, OK

$s = t$, OK

$s < t$, NOT OK.

∴ eg. $1u - 2u > 0u$

$(-1u)$?

$0x\text{FFFFFF}$ >

$0x\text{00000000}$

↓

TRUE X

Connection

```
return strlen(s) > strlen(t);
```

Don't mix unsigned int and int!
(size_t)

* Security vulnerability in getpeername.

(15) INTEGER ARITHMETIC

1. Unsigned Addition

$w=4$

0000 to 1111

(ie) 0 to 15

$\therefore \text{max sum} = 15 + 15$

30

needs 5 bits

Word size inflation

w bits + w bits \rightarrow needs $(w+1)$ bits

~~32~~ $\frac{1}{16}$ $\frac{1}{8}$ $\frac{1}{4}$ $\frac{1}{2}$ $\frac{0}{1}$

$(w+1) + (w+1) \rightarrow$ needs $(w+2)$ bits.

How to do unsigned arithmetic?

Use modular arithmetic.

eg. $\begin{array}{r} [1001] \\ + [1100] \\ \hline \end{array}$

9

12

10101

21

\rightarrow needs 5 bits

$$21 \% 16 = 5$$

2^4

\rightarrow discard the last bit \equiv subtracting 2^w from the sum.

$$x + y = \begin{cases} x+y, & x+y < 2^w \\ x+y-2^w, & 2^w \leq x+y < 2^{w+1} \end{cases}$$

\nearrow unsigned addition

\nwarrow w bits

Overflow of arithmetic operation - Full integer result cannot fit with the word size. limits of the data type.

Two's complement addition (16)

Range of bit int values: $-2^{w-1} \leq x, y \leq 2^{w-1} - 1$

Range of sum: $-2^w \leq x+y \leq 2^w - 2$
 \Rightarrow require $(w+1)$ bits to store the result.

$$\begin{aligned} & \because 2^{w-1} + 2^{w-1} \\ & = 2 \left(2^{w-1} \right) \\ & = 2^1 \cdot 2^{w-1} = 2^w \\ & = \boxed{2^w} \end{aligned}$$

$x +_w y = x+y-2^w, \quad 2^w \leq x+y$ - positive overflow
 $x +_w y = x+y, \quad -2^{w-1} \leq x+y < 2^{w-1}$ - Normal
 $x +_w y = x+y+2^w, \quad x+y < -2^{w-1}$ - Negative overflow.

eg.

x	y	$x+y$	$x +_4 y$	
$\overset{-8}{[1000]}$	$\overset{-5}{[1011]}$	$\overset{-13}{[10011]}$	$\overset{3}{[0011]}$	$\begin{array}{r} 1000 \\ 1011 \\ \hline 0011 \end{array}$
$\overset{-8}{[1000]}$	$\overset{5}{[0101]}$	$\overset{-3}{[01101]}$	$\overset{-3}{[1101]}$	$\begin{array}{r} 1000 \\ 0101 \\ \hline 1101 \end{array}$
$\overset{2}{[0010]}$	$\overset{5}{[0101]}$	$\overset{7}{[00111]}$	$\overset{7}{[0111]}$	
$\overset{5}{[0101]}$	$\overset{5}{[0101]}$	$\overset{10}{[01010]}$	$\overset{-6}{[1010]}$	$\begin{array}{r} 0101 \\ 0101 \\ \hline 1010 \end{array}$

$+3+16=3$

$10-16=-6$

(17) Two's - Complement Negation.

x 's additive inverse is $-x$

$$\therefore x + (-x) = 0$$

For all x in $-2^{w-1} < x < 2^{w-1}$, additive inv is $-x$.

But for $x = -2^{w-1}$, add. inv is -2^{w-1} .

$\therefore -2^{w-1} + -2^{w-1} = -2^w$ negative overflow so add 2^w

$$= -2^w + 2^w = 0.$$

Unsigned Multiplication

$$0 \leq x, y \leq 2^w - 1$$

$$0 \leq x \cdot y \leq (2^w - 1)^2 = 2^{2w} - 2^{w+1} + 1$$

needs $2w$ bits to represent.

Solution

Truncation of higher order bits. eg. if $w = 2$

$$x \times_w y = (x \cdot y) \bmod 2^w$$

$$2^4 - 2^3 + 1 = 16 - 8 + 1$$

$$= 9$$

need 4 bits

eg. $w = 2$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 9$$

$$= 9 \% (2^2) = 9 \% 4 = 1$$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Two's-Complement Multiplication

w=3
unsigned
2's comp.

x	y	x.y	x.y truncated
5 [101]	3 [011]	15 [001111]	7 [111]
-3 [101]	3 [011]	-9 [110111]	-1 [111]

* Multiplying by constants

x * 14

14 = 2³ + 2² + 2¹

⊗ = Integer multiplication
 eg. 1 * 14 = 14
 [0001] [1110] = [1110]

takes 10 or more CPU cycles.
 (x << 2) + (x << 1)
 [0001] << 3 = [1000]
 [0001] << 2 = [0100]
 [0001] << 1 = [0010]
 +
 [1110]

Also: 14 = 2⁴ - 2¹

so, x * 14 = (x << 4) - (x << 1)

[0001] << 4 = 10000
 [0001] << 1 = 0010 -
 [1110]

if x=2 [0010]

2 * 14 = 28 = 28 % 16 = 12
 [0010] [1110] [11100] = [1100]

= [0010] << 3 + [0010] << 2 + [0010] << 1
 = 10000 + 1000 + 0100 = 11100 = [1100]
 12.

Dividing by Powers of Two

Integer division - 30 or more clock cycles.

unsigned division by $2^k \equiv x \gg k$
(ie) $x / 2^k$ (logical right shift)

!!! for +ve signed number.

(*) Negative signed numbers (arithmetic right shift).
 $-7 / 2 = -3$

But $[1001] / 2^1 = [1001] \gg 1$ ($\because k=1$)
 $= [1100] = -4$ Wrong!

Biasing

ie. if $x < 0$, then

$$x = x + (1 \ll k) - 1$$

eg $x = -7 < 0$, $\therefore x = [1001] + ([0001] \ll 1) - [0001]$
 $= [1001] + [0010] - [0001]$
 $= [1001] + [0001] =$
 $= [1010]$

Now $x / 2^1 = [1010] \gg 1 = [1101] = -3$ Now!