



Strings and Memory Representations

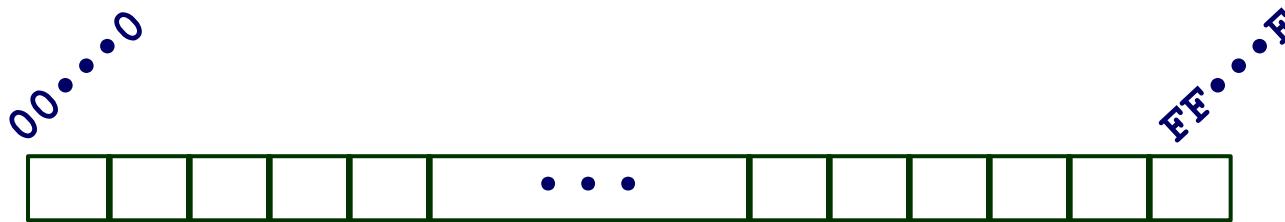
Reading: CS:APP Chapter 2.1

These slides adapted from materials provided by the textbook authors.

Strings and Memory Representations

- From Bytes to Words
- Byte ordering
- Pointer representation
- Character and string representation
- Summary

Byte-Oriented Memory Organization



- **Programs refer to data by address**
 - Conceptually, envision it as a very large array of bytes
 - In reality, it's not, but can think of it that way
 - An address is like an index into that array
 - and, a pointer variable stores an address
- **Note: system provides private address spaces to each “process”**
 - Think of a process as a program being executed
 - So, a program can clobber its own data, but not that of others

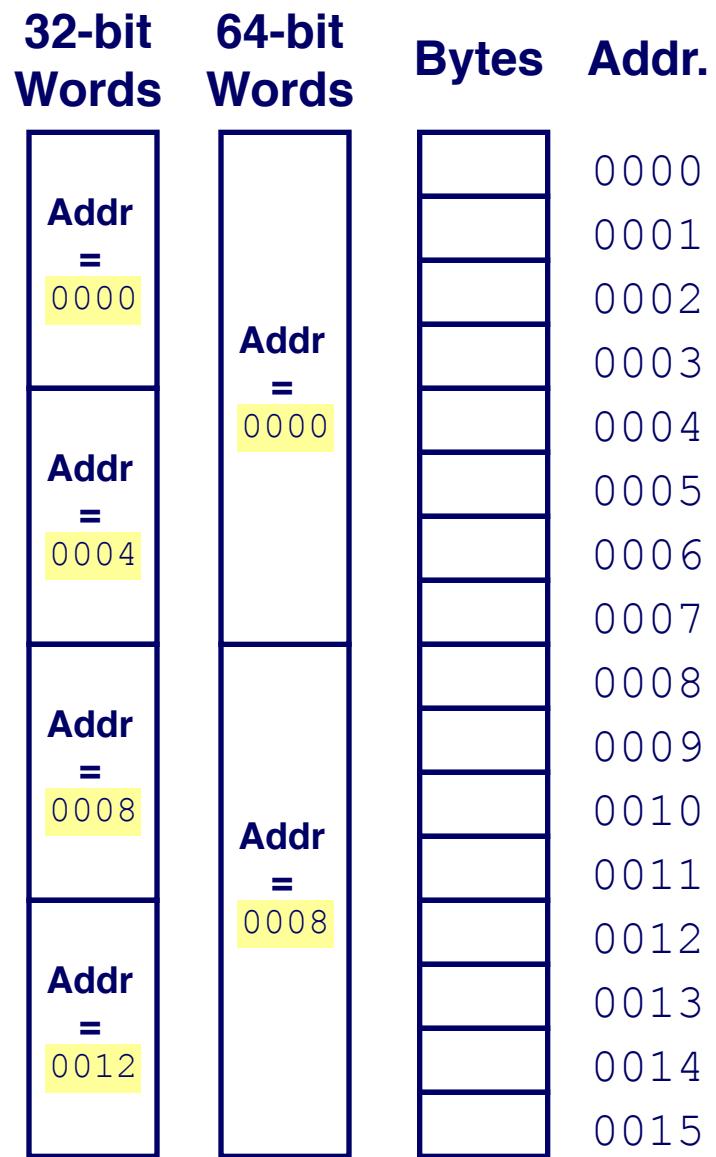
Machine Words

- Any given computer has a “Word Size”
 - Nominal size of integer-valued data
 - and of addresses
 - Until recently, most machines used 32 bits (4 bytes) as word size
 - Limits addresses to 4GB (2^{32} bytes)
 - Increasingly, machines have 64-bit word size
 - Potentially, could have 18 EB (exabytes) of addressable memory
 - That's 18.4×10^{18}
 - Machines still support multiple data formats
 - Fractions or multiples of word size
 - Always integral number of bytes

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)



Example Data Representations

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

Byte Ordering

- So, how are the bytes within a multi-byte word ordered in memory?
- Conventions
 - Big Endian: Sun, PPC Mac, Internet
 - Least significant byte has highest address
 - Little Endian: x86, ARM processors running Android, iOS, and Windows
 - Least significant byte has lowest address

Byte Ordering Example

■ Example

- Variable Y has 4-byte value of 0x01234567
- Address given by &Y is 0x100



BigEndian

0x100 0x101 0x102 0x103



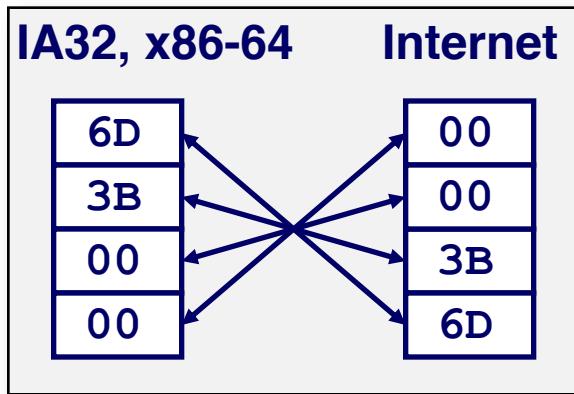
LittleEndian

0x100 0x101 0x102 0x103

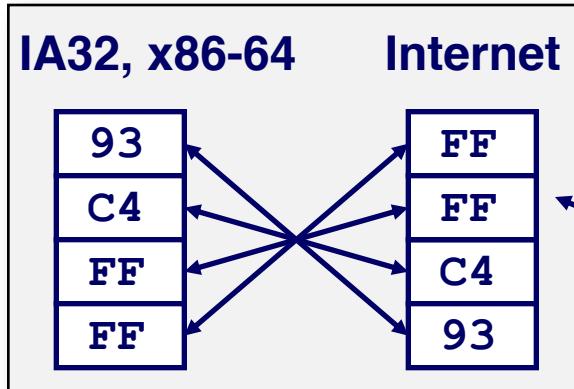


Representing Integers

```
int A = 15213;
```



```
int B = -15213;
```

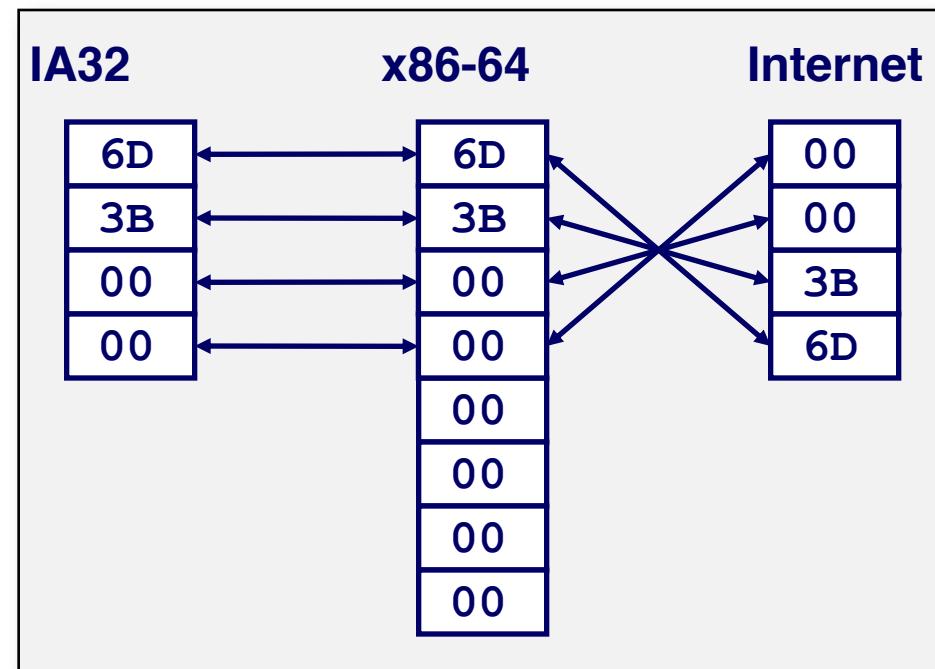


Decimal: 15213

Binary: 0011 1011 0110 1101

Hex: 3 B 6 D

```
long int C = 15213;
```



Two's complement representation

Examining Data Representations

■ Code to Print Byte Representation of Data

- Casting pointer to unsigned char * allows treatment as a byte array

```
typedef unsigned char *pointer;

void show_bytes(pointer start, size_t len) {
    size_t i;
    for (i = 0; i < len; i++)
        printf("%p\t0x%.2x\n", start+i, start[i]);
    printf("\n");
}
```

Printf directives:

%p: Print pointer
%x: Print Hexadecimal

show_bytes Execution Example

```
int a = 15213;  
printf("int a = 15213;\\n");  
show_bytes((pointer) &a, sizeof(int));
```

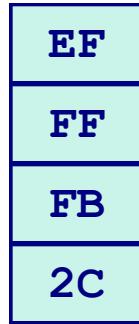
Result (Linux x86-64):

```
int a = 15213;  
0x7ffb7f71dbc      6d  
0x7ffb7f71dbd      3b  
0x7ffb7f71dbe      00  
0x7ffb7f71dbf      00
```

Representing Pointers

```
int B = -15213;  
int *P = &B;
```

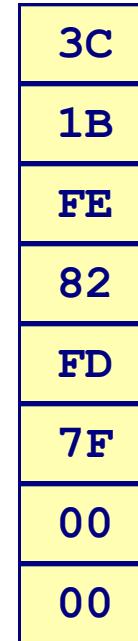
Internet / Sun



IA32



x86-64



Different compilers & machines assign different locations to objects

Even get different results each time run program

Representing Characters

- Representing letters and words by numbers: morse code
- Later TTY systems converted typewriter into morse / baudot other codes
- Led to ASCII (American Standard Code for Information Interchange) in ~1960
- 7 (then 8) bit code for letters, numbers, “actions”
- Expanded to Unicode in 90’s and UTF-8 in ‘00’s

ଫ୍ରେଣ୍ଟ ଏଣ୍ଡ ଏମ୍ୟ



USASCII code chart							
b ₇ b ₆ b ₅			b ₄	b ₃	b ₂	b ₁	Column Row
0	0	0	0	0	1	0	1
0	0	0	0	0	NUL	DLE	SP
0	0	0	1	1	SOH	DC1	!
0	0	1	0	2	STX	DC2	"
0	0	1	1	3	ETX	DC3	#
0	1	0	0	4	EOT	DC4	\$
0	1	0	1	5	ENQ	NAK	%
0	1	1	0	6	ACK	SYN	&
0	1	1	1	7	BEL	ETB	'
1	0	0	0	8	BS	CAN	(
1	0	0	1	9	HT	EM)
1	0	1	0	10	LF	SUB	*
1	0	1	1	11	VT	ESC	:
1	1	0	0	12	FF	FS	,
1	1	0	1	13	CR	GS	-
1	1	1	0	14	SO	RS	>
1	1	1	1	15	SI	US	?

ASCII Table == First 127 of UTF-8

	Alphabetic		Extended punctuation
	Control character		Graphic character
	Numeric digit		International
	Punctuation		Undefined

ASCII (1977/1986)																
	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
0_	NUL 0000	SOH 0001	STX 0002	ETX 0003	EOT 0004	ENQ 0005	ACK 0006	BEL 0007	BS 0008	HT 0009	LF 000A	VT 000B	FF 000C	CR 000D	SO 000E	SI 000F
	0 1	1 2	2 3	3 4	4 5	5 6	6 7	7 8	8 9	9 10	10 11	11 12	12 13	13 14	14 15	
	DLE 0010	DC1 0011	DC2 0012	DC3 0013	DC4 0014	NAK 0015	SYN 0016	ETB 0017	CAN 0018	EM 0019	SUB 001A	ESC 001B	FS 001C	GS 001D	RS 001E	US 001F
1_	16 17	17 18	18 19	19 20	20 21	21 22	22 23	23 24	24 25	25 26	26 27	27 28	28 29	29 30	30 31	
2_	SP 0020	! 0021	" 0022	# 0023	\$ 0024	% 0025	& 0026	' 0027	(0028) 0029	* 002A	+002B	,002C	- 002D	. 002E	/ 002F
3_	32 33	33 34	34 35	35 36	36 37	37 38	38 39	39 40	40 41	41 42	42 43	43 44	44 45	45 46	46 47	
4_	0 0030	1 0031	2 0032	3 0033	4 0034	5 0035	6 0036	7 0037	8 0038	9 0039	:	; 003A	< 003B	= 003C	> 003D	? 003E
5_	48 49	49 50	50 51	51 52	52 53	53 54	54 55	55 56	56 57	57 58	58 59	59 60	60 61	61 62	62 63	
6_	@ 0040	A 0041	B 0042	C 0043	D 0044	E 0045	F 0046	G 0047	H 0048	I 0049	J 004A	K 004B	L 004C	M 004D	N 004E	O 004F
7_	64 65	65 66	66 67	67 68	68 69	69 70	70 71	71 72	72 73	73 74	74 75	75 76	76 77	77 78	78 79	
8_	P 0050	Q 0051	R 0052	S 0053	T 0054	U 0055	V 0056	W 0057	X 0058	Y 0059	Z 005A	[005B	\ 005C] 005D	^ 005E	— 005F
9_	80 81	81 82	82 83	83 84	84 85	85 86	86 87	87 88	88 89	89 90	90 91	91 92	92 93	93 94	94 95	
10_	~ 0060	a 0061	b 0062	c 0063	d 0064	e 0065	f 0066	g 0067	h 0068	i 0069	j 006A	k 006B	l 006C	m 006D	n 006E	o 006F
11_	96 97	97 98	98 99	99 100	100 101	101 102	102 103	103 104	104 105	105 106	106 107	107 108	108 109	109 110	110 111	
12_	p 0070	q 0071	r 0072	s 0073	t 0074	u 0075	v 0076	w 0077	x 0078	y 0079	z 007A	{ 007B	 007C	} 007D	~ 007E	DEL 007F
13_	112 113	113 114	114 115	115 116	116 117	117 118	118 119	119 120	120 121	121 122	122 123	123 124	124 125	125 126	126 127	

ASCII collating order

- Sequences 0..9, a..z and A..Z are in order

- Checking if 'd' is a digit:

`(d >= '0') && (d <= '9')`

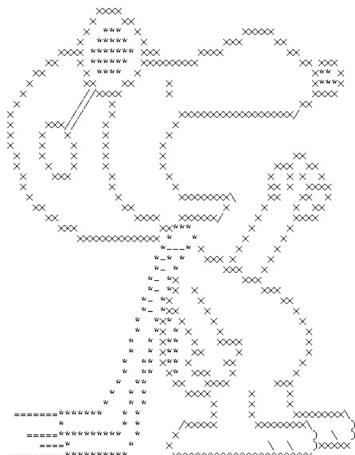
- Checking if 'l' is a letter:

`((l >= 'a') && (l <= 'z')) || ((l >= 'A') && (l <= 'Z'))`

- Change letter 'A'..'Z' to lower case:

`(l - 'A') + 'a'`

Humans make art no matter how primitive the tools



<http://patorjk.com/software/taag/#p=display&f=Doh&t=csci%202400>

Representing Strings

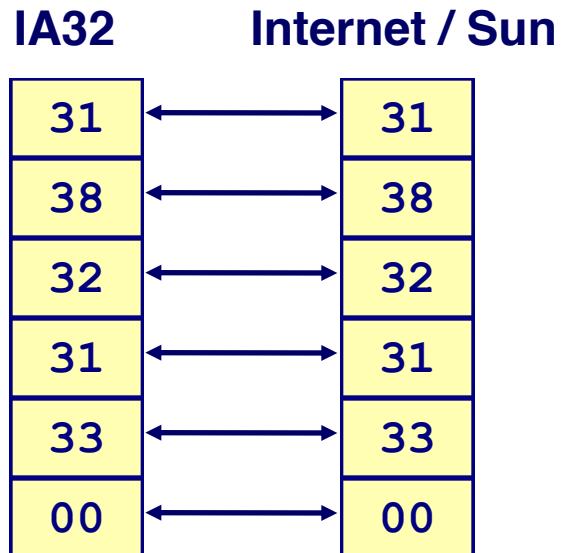
```
char S[6] = "18213";
```

■ Strings in C

- Represented by array of 'char'
- Char 'x' enclosed in single quotes, string "enclosed in double"
- 'abcd' is not same as "abcd"
 - 'abcd' is a 4-byte integer, not a string
 - "abcd" is 5 byte string with null terminator
- Each character encoded in ASCII format
- String should be null-terminated
 - Final character = 0

■ Compatibility

- Byte ordering not an issue for strings



C sizeof() vs. string length vs. constant

- Common mistake is to think sizeof() returns the string length or the length of the “buffer” that is pointed to
- A “char *” points to a char; convention is it’s the first char of a string but C doesn’t keep track of length

```
int main()
{
    char buffer[100];
    char *ptr;
    char *str = "how";
    int istr = 'how';
    printf("sizeof(buffer) = %d\n", sizeof(buffer));
    printf("sizeof(ptr)    = %d\n", sizeof(ptr));
    printf("sizeof(*ptr)   = %d\n", sizeof(*ptr));
    printf("sizeof(str)    = %d\n", sizeof(str));
    printf("sizeof(istr)   = %d\n", sizeof(istr));
}
```

String and Memory Functions in C

<http://www.cplusplus.com/reference/cstring/>

String

- **strcpy(dst, src)** – copy
- **strlen(str)** – length
- **strcat(dest, src)** – append
- **strcmp(src1, src2)** – compare

Memory

- **memcpy(dst, src, len)**
- **memmove(dst, src, len)**
- **memset(dst, value, len)**
- **memcmp(dst, src, len)**

- **strnlen(str, max)**
- **strncpy(dst, src, dstsize)**
- **strncat(dest, src, dstsize)**
- **strncmp(src1, src2, max)**

C string functions – strcpy 3 ways

■ And all of them bad!

```
char * strcpy(char *dst, char *src) {
    int i;
    for( i = 0; src[i] != '\0'; i++) {
        dst[i] = src[i];
    }
    dst[i] = 0;
    return dst;
}
```

```
char * strcpy(char *dst, char *src) {
    char *orig = dst;
    for( ; *src ; dst++, src++ ) {
        *dst = *src;
    }
    *dst = 0;
    return orig;
}
```

```
char * strcpy(char *dst, char *src) {
    char *orig = dst;
    while(*src) {
        *dst++ = *src++;
    }
    *dst = 0;
    return orig;
}
```

C string functions – strcpy

- Most compilers will warn about strcpy, strlen, etc

```
char * strcpy(char *dst, char *src, size_t len) {
    int i;

    for( i = 0; i < len && src[i] != '\0'; i++) {
        dst[i] = src[i];
    }
    dst[i] = 0;
    return dst;
}
```

- You'll be directed to the 'n' versions

C memory functions – memcpy

- The mem* functions don't look for null terminator

```
void* memcpy(void *dst, void *src, size_t len) {
    int i;
    char *cdst = (char *) dst;
    char *csrc = (char *) src;
    for( i = 0; i < len ; i++) {
        cdst[i] = csrc[i];
    }
    return dst;
}
```

- Use provided str* and mem* functions rather than writing your own.
 - Much, much faster code
 - They already made it correct

C strcmp and memcmp

- **strncmp(a,b) and memcmp(a,b,len) return -1, 0 or 1**
 - -1 a < b
 - 0 0 == b
 - 1 a > b
- Leads to weird looking programming idiom

```
if ( ! strncmp(a,b,n) ) {  
    /* a is equal to b */  
} else {  
    /* a != b */  
}
```

- Better to write this like..

```
if ( ! strncmp(a,b,n)== 0 ) {  
    /* a is equal to b */  
}
```

C memcpy and memmove

- **memcpy(a,b,len) and memmove(a,b,len) are similar**
 - memcpy will blindly copy ‘b’ to ‘a’ up to ‘len’ bytes
 - memmove acts as if data was copied to intermediate location.
This is important if ‘a’ and ‘b’ overlap
- **Example: copying a buffer to make room in first element**

```
int buffer[N] ;  
....  
memmove(&buffer[1], &buffer[0], (N-1) * sizeof(int))
```

C `strncat`

- `strncat(a,b,len)` appends ‘b’ to ‘a’
 - Must have enough room in ‘a’
- You can accomplish same goal using `snprintf()`

```
char buffer[1024];
snprintf(a, len, "%s%s", a, b)
```