# Bits, Bytes, and Integers

# Today: Bits, Bytes, and Integers

- **Representing information as bits**
- **Bit-level manipulations**
- **Integers**
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
- **Summary**

# Number representations

- Understand the ranges of values that can be represented and the properties of the different arithmetic operations.

- This understanding is critical to writing programs that work correctly over the full range of numeric values and that are portable across different combinations of machine, operating system, and compiler

- Whereas in an earlier era program bugs would only inconvenience people when they happened to be triggered, there are now legions of hackers who try to exploit any bug they can find to obtain unauthorized access to other people's systems.

- This puts a higher level of obligation on programmers to understand how their programs work and how they can be made to behave in undesirable ways.
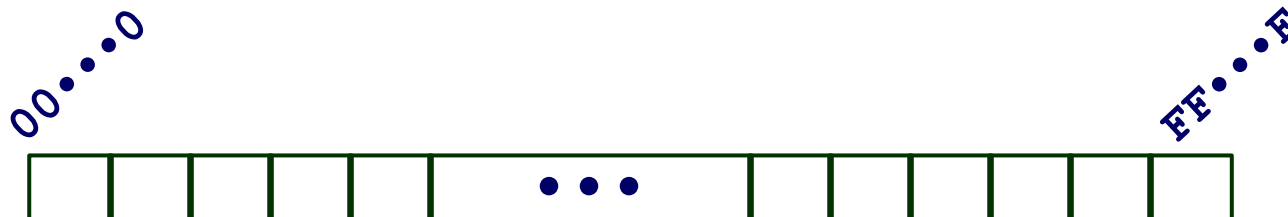
# Encoding Byte Values

- **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 |

# 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
  - Machines support multiple data formats
    - Fractions or multiples of word size
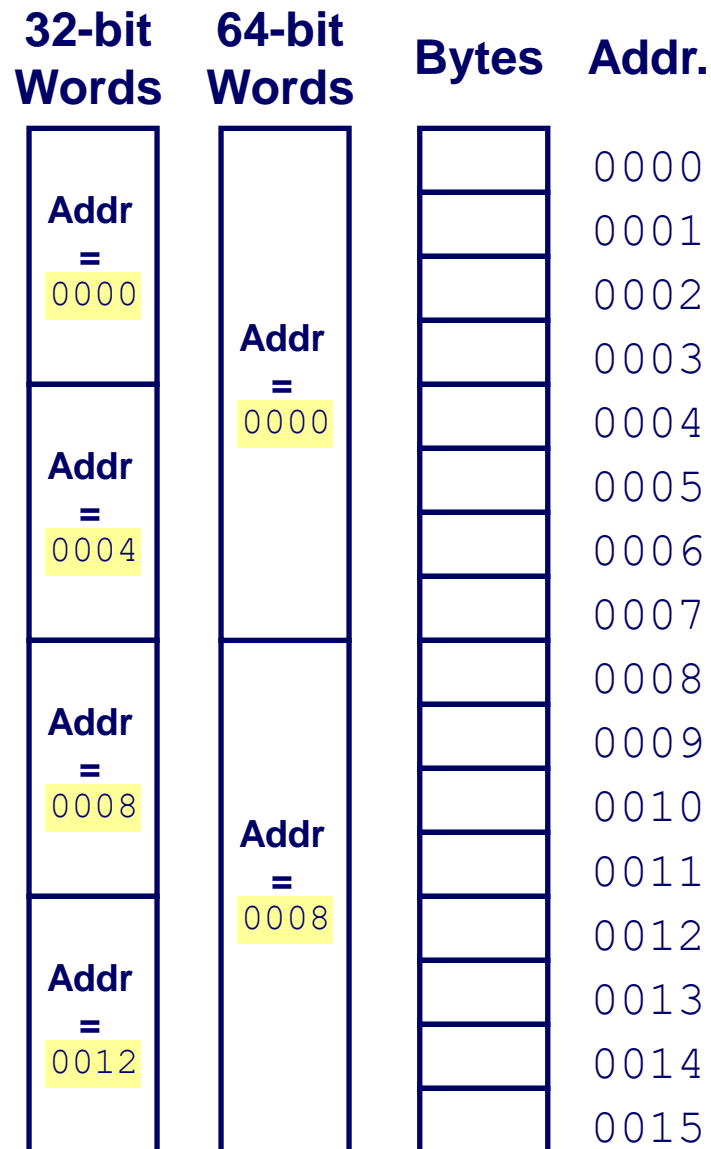    - Always integral number of bytes

# Address accessible and Word size

- **32 bit word size**

- **Address range: 0-2^32-1**

- **4 GB RAM**


- **Check your PCs specifications**

# 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 |
| **Addr =** 0008 | | | 0009 |
| | | | 0010 |
| | **Addr =** 0008 | | 0011 |
| | | | 0012 |
| **Addr =** 0012 | | | 0013 |
| | | | 0014 |
| | | | 0015 |

8

In recent years, there has been a widespread shift from machines with 32-bit word sizes to those with word sizes of 64 bits.

Most 64-bit machines can also run programs compiled for use on 32-bit machines, a form of backward compatibility. So, for example, when a program `prog.c` is compiled with the directive

> linux> *gcc -m32 prog.c*

then this program will run correctly on either a 32-bit or a 64-bit machine. On the other hand, a program compiled with the directive

> linux> *gcc -m64 prog.c*

will only run on a 64-bit machine. We will therefore refer to programs as being either "32-bit programs" or "64-bit programs," since the distinction lies in how a program is compiled, rather than the type of machine on which it runs.

# 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 |

# Byte Ordering

- **How should bytes within a multi-byte word be ordered in memory?**

- **Example:** suppose a variable x of type int has address 0x100, that is, the value of the address expression &x is 0x100. Then the 4 bytes of x would be stored in memory locations 0x100, 0x101, 0x102, and 0x103.

- **Conventions**

  - Big Endian: Sun, PPC Mac, Internet

    - most significant byte comes first

  - Little Endian: x86

    - least significant byte comes first

the terms "little endian" and "big endian" come from the book *Gulliver's Travels* by Jonathan Swift

SIMPLY EXPLAINED

BIG-ENDIAN                LITTLE-ENDIAN

0xCAFEBABE                0xCAFEBABE
will be stored as         will be stored as
CA | FE | BA | BE         BE | BA | FE | CA

www.thebittheories.com

# Byte Ordering Example

- **Big Endian**
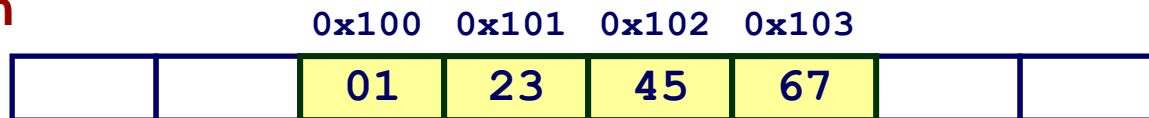  - Least significant byte has highest address

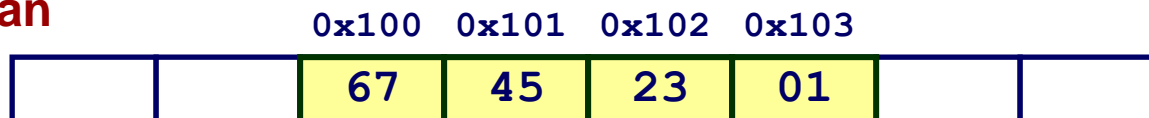- **Little Endian**
  - Least significant byte has lowest address

- **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 | | |

At times, byte ordering becomes an issue. The first is when binary data are communicated over a network between different machines.

A common problem is for data produced by a little-endian machine to be sent to a big-endian machine, or vice versa, leading to the bytes within the words being in reverse order for the receiving program.

To avoid such problems, code written for networking applications must follow established conventions for byte ordering to make sure the sending machine converts its internal representation to the network standard, while the receiving machine converts the network standard to its internal representation

# Reading Byte-Reversed Listings

**Objdump –d  test.o**

- ## 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 00 | cmpl   $0x0,0x28(%ebx) |

- ## Deciphering Numbers
  - Value:                                          0x12ab
  - Pad to 32 bits:                    0x000012ab
  - Split into bytes:                   00 00 12 ab
  - Reverse:                               ab 12 00 00

# Examining Data Representations

■ **Code to Print Byte Representation of Data**

  ▪ Casting pointer to unsigned char * creates byte array

```
typedef unsigned char *pointer;

void show_bytes(pointer start, int len){
  int 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):**

```
int a = 15213;
0x11ffffcb8 0x6d
0x11ffffcb9 0x3b
0x11ffffcba 0x00
0x11ffffcbb 0x00
```

```
1    void test_show_bytes(int val) {
2        int ival = val;
3        float fval = (float) ival;
4        int *pval = &ival;
5        show_int(ival);
6        show_float(fval);
7        show_pointer(pval);
8    }
```

**What do you notice?**

| Machine | Value | Type | Bytes (hex) |
|---------|-------|------|-------------|
| Linux 32 | 12,345 | int | 39 30 00 00 |
| Windows | 12,345 | int | 39 30 00 00 |
| Sun | 12,345 | int | 00 00 30 39 |
| Linux 64 | 12,345 | int | 39 30 00 00 |
| Linux 32 | 12,345.0 | float | 00 e4 40 46 |
| Windows | 12,345.0 | float | 00 e4 40 46 |
| Sun | 12,345.0 | float | 46 40 e4 00 |
| Linux 64 | 12,345.0 | float | 00 e4 40 46 |
| Linux 32 | &ival | int * | e4 f9 ff bf |
| Windows | &ival | int * | b4 cc 22 00 |
| Sun | &ival | int * | ef ff fa 0c |
| Linux 64 | &ival | int * | b8 11 e5 ff ff 7f 00 00 |

# 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 | |

int B = -15213;

**IA32, x86-64**     **Sun**

| 93 | | FF |
|---|---|---|
| C4 | | FF |
| FF | | C4 |
| FF | | 93 |

**Two's complement representation (Covered later)**

# Representing Pointers

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

| Sun | IA32 | x86-64 |
|:---:|:---:|:---:|
| EF | D4 | 0C |
| FF | F8 | 89 |
| FB | FF | EC |
| 2C | BF | FF |
| | | FF |
| | | 7F |
| | | 00 |
| | | 00 |

**Different compilers & machines assign different locations to objects**
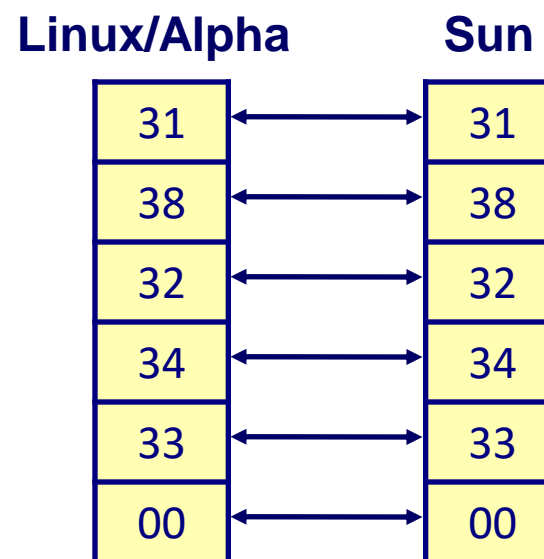
# Representing Strings

```
char S[6] = "18243";
```

- **Strings in C**
  - Represented by array of characters
  - Each character encoded in ASCII format
    - 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**

| Linux/Alpha | | Sun |
|:---:|:---:|:---:|
| 31 | ⟷ | 31 |
| 38 | ⟷ | 38 |
| 32 | ⟷ | 32 |
| 34 | ⟷ | 34 |
| 33 | ⟷ | 33 |
| 00 | ⟷ | 00 |

If you have a simple 8-bit character representation (e.g. extended ASCII), then no, endianness does not affect the layout, because each character is one byte.
If you have a multi-byte representation, such as UTF-16, then yes, endianness is still important

# Today: Bits, Bytes, and Integers

- **Representing information as bits**
- **Bit-level manipulations**
- **Integers**
  - Representation: unsigned and signed
  - Conversion, casting
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
- **Summary**

# Boolean Algebra

- **Developed by George Boole in 19th Century, applied to logic reasoning**
  - 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 | 1 |

**Or**

- **A|B = 1 when either A=1 or B=1**

| \| | 0 | 1 |
|---|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 1 |

**Not**

- **~A = 1 when A=0**

| ~ | |
|---|---|
| 0 | 1 |
| 1 | 0 |

**Exclusive-Or (Xor)**

- **A^B = 1 when either A=1 or B=1, but not both**

| ^ | 0 | 1 |
|---|---|---|
| 0 | 0 | 1 |
| 1 | 1 | 0 |

# Application of Boolean Algebra

- **Applied to Digital Systems by Claude Shannon**
  - 1937 MIT Master's Thesis
  - Reason about networks of relay switches
    - Encode closed switch as 1, open switch as 0

**A&~B**

A&~B

~A&B

~B

B

**~A&B**

**Connection when**

**A&~B | ~A&B**

**= A^B**

Carnegie Mellon

# General Boolean Algebras

■ **Operate on Bit Vectors**

▪ Operations applied bitwise

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

■ **All of the Properties of Boolean Algebra Apply**

# 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 (Char data type)**
  - ~0x41 = 0xBE
    - ~$01000001_2$ = $10111110_2$
  - ~0x00 = 0xFF
    - ~$00000000_2$ = $11111111_2$
  - 0x69 & 0x55 = 0x41
    - $01101001_2$ & $01010101_2$ = $01000001_2$
  - 0x69 | 0x55 = 0x7D
    - $01101001_2$ | $01010101_2$ = $01111101_2$

# Contrast: Logic Operations in C

- **Contrast to Logical Operators**
  - &&, ||, !
    - View 0 as "False"
    - Anything nonzero as "True"
    - Always return 0 or 1

- **Examples (char data type)**
  - !0x41 = 0x00
  - !0x00 = 0x01
  - !!0x41 = 0x01

  - 0x69 && 0x55 = 0x01
  - 0x69 || 0x55 = 0x01
  - p && *p        (avoids null pointer access) <span style="color:red">(short-circuit)</span>

# 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

- **Undefined Behavior**
  - Shift amount < 0 or ≥ word size

| Argument x | 01100010 |
|---|---|
| << 3 | 00010*000* |
| Log. >> 2 | *00*011000 |
| Arith. >> 2 | *00*011000 |

| Argument x | 10100010 |
|---|---|
| << 3 | 00010*000* |
| Log. >> 2 | *00*101000 |
| Arith. >> 2 | *11*101000 |

# Today: Bits, Bytes, and Integers

- **Representing information as bits**
- **Bit-level manipulations**
- **Integers**
    - **Representation: unsigned and signed**
    - Conversion, casting
    - Expanding, truncating
    - Addition, negation, multiplication, shifting
- **Summary**

# Encoding Integers

**Unsigned**

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

**Two's Complement**

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

**Sign Bit**

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

- **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

# Encoding Example (Cont.)

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

| Weight | 15213 | | -15213 | |
|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 |
| 2 | 0 | 0 | 1 | 2 |
| 4 | 1 | 4 | 0 | 0 |
| 8 | 1 | 8 | 0 | 0 |
| 16 | 0 | 0 | 1 | 16 |
| 32 | 1 | 32 | 0 | 0 |
| 64 | 1 | 64 | 0 | 0 |
| 128 | 0 | 0 | 1 | 128 |
| 256 | 1 | 256 | 0 | 0 |
| 512 | 1 | 512 | 0 | 0 |
| 1024 | 0 | 0 | 1 | 1024 |
| 2048 | 1 | 2048 | 0 | 0 |
| 4096 | 1 | 4096 | 0 | 0 |
| 8192 | 1 | 8192 | 0 | 0 |
| 16384 | 0 | 0 | 1 | 16384 |
| -32768 | 0 | 0 | 1 | -32768 |
| **Sum** | | **15213** | | **-15213** |

| C data type | Minimum | Maximum |
|---|---|---|
| char | −128 | 127 |
| unsigned char | 0 | 255 |
| short [int] | −32,768 | 32,767 |
| unsigned short [int] | 0 | 65,535 |
| int | −2,147,483,648 | 2,147,483,647 |
| unsigned [int] | 0 | 4,294,967,295 |
| long [int] | −2,147,483,648 | 2,147,483,647 |
| unsigned long [int] | 0 | 4,294,967,295 |
| long long [int] | −9,223,372,036,854,775,808 | 9,223,372,036,854,775,807 |
| unsigned long long [int] | 0 | 18,446,744,073,709,551,615 |

Figure 2.8 **Typical ranges for C integral data types on a 32-bit machine.** Text in square brackets is optional.

# 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

- **Other Values**
  - Minus 1
    
    111...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 |

# 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

# 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

# Today: Bits, Bytes, and Integers

- **Representing information as bits**
- **Bit-level manipulations**
- **Integers**
  - Representation: unsigned and signed
  - **Conversion, casting**
  - Expanding, truncating
  - Addition, negation, multiplication, shifting
- **Summary**

# Signed numbers

- **Computer programs calculate both positive and negative numbers, so we need a representation that distinguishes the positive from the negative.**

- **Add a separate sign, which conveniently can be represented in a single bit; the name for this representation is *sign and magnitude*.**

  - where to put the sign bit. To the right? To the left ?

  - adders for sign and magnitude may need an extra step to set the sign because we can't know in advance what the proper sign will be.

  - separate sign bit means that sign and magnitude has both a positive and a negative zero

# 2s-Complement Signed Integers

- **Given an n-bit number**

$$x = -x_{n-1}2^{n-1} + x_{n-2}2^{n-2} + \cdots + x_1 2^1 + x_0 2^0$$

- **Range: $-2^{n-1}$ to $+2^{n-1} - 1$**

- **Example**
  - **$1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1100_2$**
    **$= -1 \times 2^{31} + 1 \times 2^{30} + \ldots + 1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$**
    **$= -2{,}147{,}483{,}648 + 2{,}147{,}483{,}644 = -4_{10}$**

- **Using 32 bits**
  - **$-2{,}147{,}483{,}648$ to $+2{,}147{,}483{,}647$**

**leading 0s mean positive, and leading 1s mean negative**

# Conversion Visualized

- **2's Comp. $\rightarrow$ Unsigned**
  - Ordering Inversion
  - Negative $\rightarrow$ Big Positive



$UMax$
$UMax - 1$

$TMax + 1$
$TMax$

Unsigned Range

2's Complement Range

$TMax$

0
−1
−2

$TMin$

0

# 2s-Complement Signed Integers

- **Bit 31 is sign bit**
  - 1 for negative numbers
  - 0 for non-negative numbers

- **Non-negative numbers have the same unsigned and 2s-complement representation**

- **Some specific numbers**
  - 0:   0000 0000 … 0000
  - −1:   1111 1111 … 1111
  - Most-negative:   1000 0000 … 0000
  - Most-positive:   0111 1111 … 1111

**EXAMPLE**

**ANSWER**

**Binary to Decimal Conversion**

What is the decimal value of this 32-bit two's complement number?

$$1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1100_{two}$$

Substituting the number's bit values into the formula above:

$$(1 \times -2^{31}) + (1 \times 2^{30}) + (1 \times 2^{29}) + \ldots + (1 \times 2^{1}) + (0 \times 2^{1}) + (0 \times 2^{0})$$
$$= -2^{31} + 2^{30} + 2^{29} + \ldots + 2^{2} + 0 + 0$$
$$= -2,147,483,648_{ten} + 2,147,483,644_{ten}$$
$$= -4_{ten}$$

# Signed Negation - shortcut

- **Complement and add 1**
  - Complement means $1 \rightarrow 0$, $0 \rightarrow 1$

$$x + \overline{x} = 1111...111_2 = -1$$

$$\overline{x} + 1 = -x$$

- **Example: negate +2**
  - **+2 = 0000 0000 … 0010$_2$**
  - **−2 = 1111 1111 … 1101$_2$ + 1**
    **= 1111 1111 … 1110$_2$**

# Signed vs. Unsigned in C

- **Constants**
  - By default are considered to be **signed** integers
  - Unsigned if have "U" as suffix

    `0U, 4294967259U`

- **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;
    ```

```
1        short    int    v  = -12345;
2        unsigned short uv = (unsigned short) v;
3        printf("v = %d, uv = %u\n", v, uv);
```

When run on a two's-complement machine, it generates the following output:

```
v = -12345, uv = 53191
```

```
1            int x = -1;
2            unsigned u = 2147483648; /* 2 to the 31st */
3
4            printf("x = %u = %d\n", x, x);
5            printf("u = %u = %d\n", u, u);
```

When run on a 32-bit machine, it prints the following:

```
x = 4294967295 = -1
u = 2147483648 = -2147483648
```

In both cases, `printf` prints the word first as if it represented an unsigned number, and second as if it represented a signed number. We can see the conversion routines in action: $T2U_{32}(-1) = UMax_{32} = 2^{32} - 1$ and $U2T_{32}(2^{31}) = 2^{31} - 2^{32} = -2^{31} = TMin_{32}$.

# Casting Surprises

When an operation is performed where one operand is signed and the other is unsigned, C implicitly casts the **signed argument to unsigned** and performs the operations assuming the numbers are nonnegative.

| Expression | Type | Evaluation |
|---|---|---|
| 0 == 0U | unsigned | 1 |
| −1 < 0 | signed | 1 |
| −1 < 0U | unsigned | 0 * |
| 2147483647 > −2147483647−1 | signed | 1 |
| 2147483647U > −2147483647−1 | unsigned | 0 * |
| 2147483647 > (int) 2147483648U | signed | 1 * |
| −1 > −2 | signed | 1 |
| (unsigned) −1 > −2 | unsigned | 1 |

Figure 2.18 **Effects of C promotion rules.** Nonintuitive cases marked by '*'.

# Code Security Example

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

/* Declaration of library function memcpy */
7 void *memcpy(void *dest, void *src, size_t n);

/* 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;
}
```

- Similar to code found in FreeBSD's implementation of getpeername

- **There are legions of smart people trying to find vulnerabilities in programs**

# Typical 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);
    printf("%s\n", mybuf);
}
```

# Malicious Usage

```
/* Declaration of library function memcpy */
void *memcpy(void *dest, void *src, size_t n);
```

```
/* 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);
    . . .
}
```

# Summary
# 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`!!