# CS 240 Stage 2
# **Hardware-Software Interface**

Memory addressing, C language, pointers

Assertions, debugging

Machine code, assembly language, program translation

Control flow

Procedures, stacks

Data layout, security, linking and loading

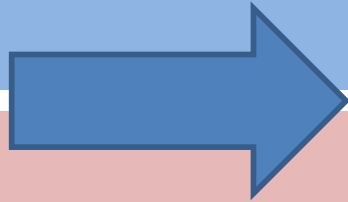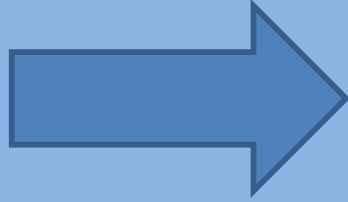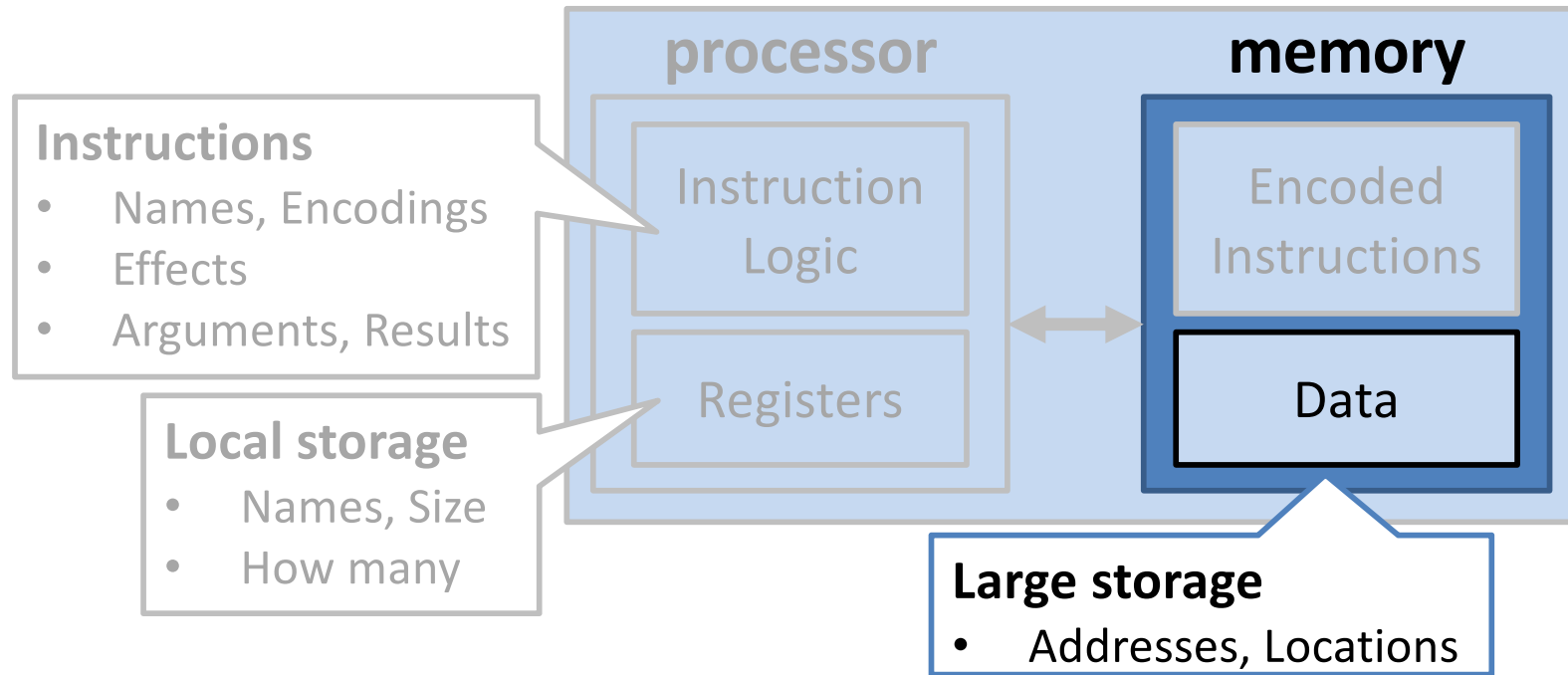| | |
|---|---|
| **Software** | Program, Application |
| | **Programming Language** |
| | **Compiler/Interpreter** |
| | Operating System |
| **Hardware** | **Instruction Set Architecture** |
| | Microarchitecture |
| | Digital Logic |
| | Devices (transistors, etc.) |
| | Solid-State Physics |

# Programming with Memory
 via C, pointers, and arrays

Why not just registers?

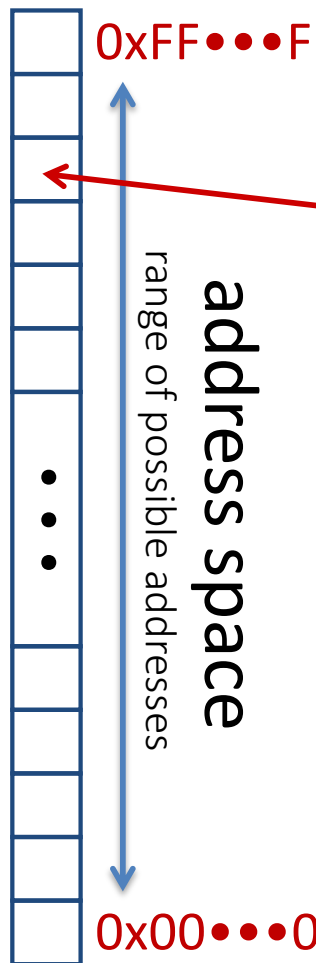- Represent larger structures
- Computable addressing
- Indirection

# Instruction Set Architecture (HW/SW Interface)

**processor**

**memory**

**Instructions**
- Names, Encodings
- Effects
- Arguments, Results

Instruction Logic

Encoded Instructions

**Local storage**
- Names, Size
- How many

Registers

Data

**Large storage**
- Addresses, Locations

# Computer

# byte-addressable memory = mutable byte array

0xFF•••F

address space
range of possible addresses

0x00•••0

**Cell** / location = element

- Addressed by unique numerical address
- Holds one byte
- Readable and writable

**Address** = index

- Unsigned number
- Represented by one word
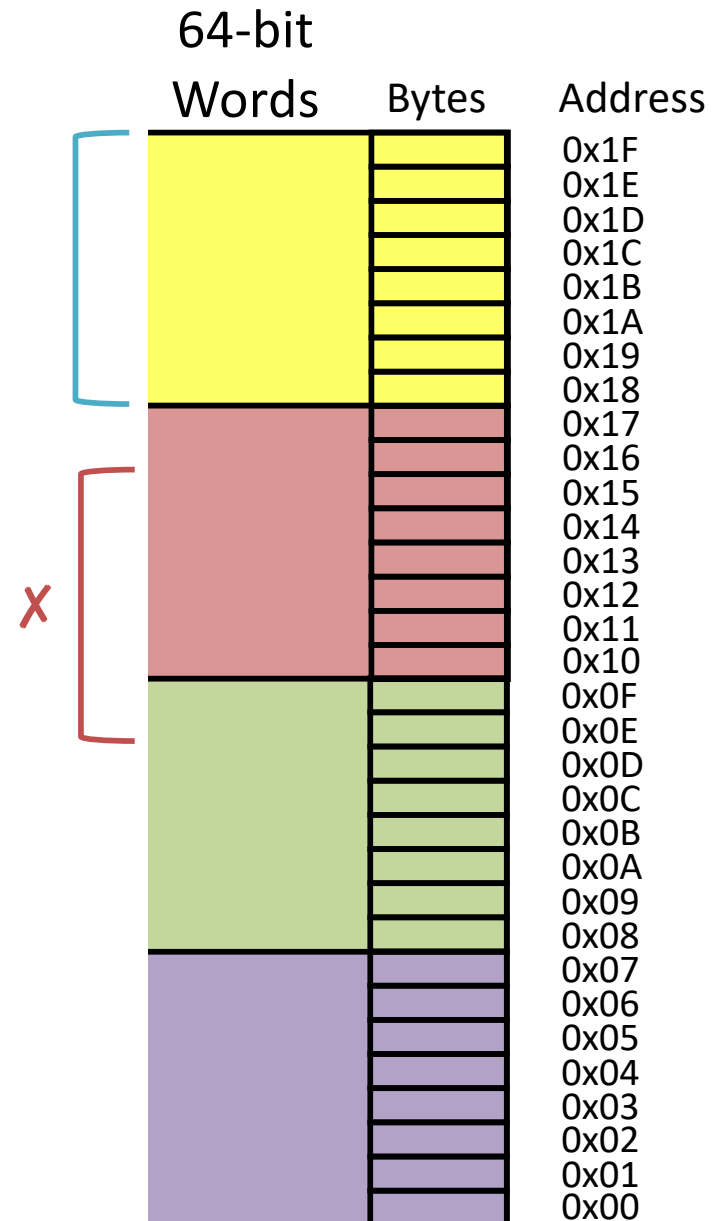- Computable and storable as a value

# multi-byte values in memory
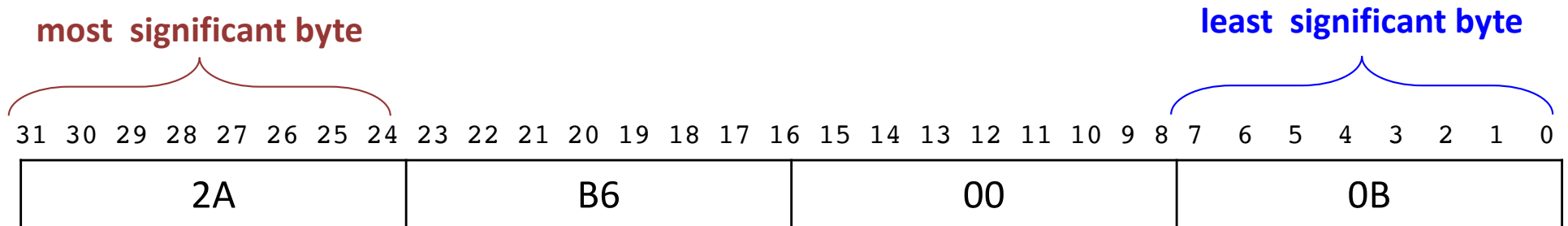
Store across contiguous byte locations.

Alignment  (Why?)

Bit order within byte always same.

Byte ordering within larger value?

64-bit Words | Bytes | Address

| Address |
|---------|
| 0x1F |
| 0x1E |
| 0x1D |
| 0x1C |
| 0x1B |
| 0x1A |
| 0x19 |
| 0x18 |
| 0x17 |
| 0x16 |
| 0x15 |
| 0x14 |
| 0x13 |
| 0x12 |
| 0x11 |
| 0x10 |
| 0x0F |
| 0x0E |
| 0x0D |
| 0x0C |
| 0x0B |
| 0x0A |
| 0x09 |
| 0x08 |
| 0x07 |
| 0x06 |
| 0x05 |
| 0x04 |
| 0x03 |
| 0x02 |
| 0x01 |
| 0x00 |

X

# *Endianness:* To store a multi-byte value in memory, which byte is stored first (at a lower address)?

most  significant byte                                                                  least  significant byte

| 31 30 29 28 27 26 25 24 | 23 22 21 20 19 18 17 16 | 15 14 13 12 11 10 9 8 | 7 6 5 4 3 2 1 0 |
|---|---|---|---|
| 2A | B6 | 00 | 0B |

| Address | Contents |
|---|---|
| 03 | 2A |
| 02 | B6 |
| 01 | 00 |
| 00 | 0B |

| Address | Contents |
|---|---|
| 03 | 0B |
| 02 | 00 |
| 01 | B6 |
| 00 | 2A |

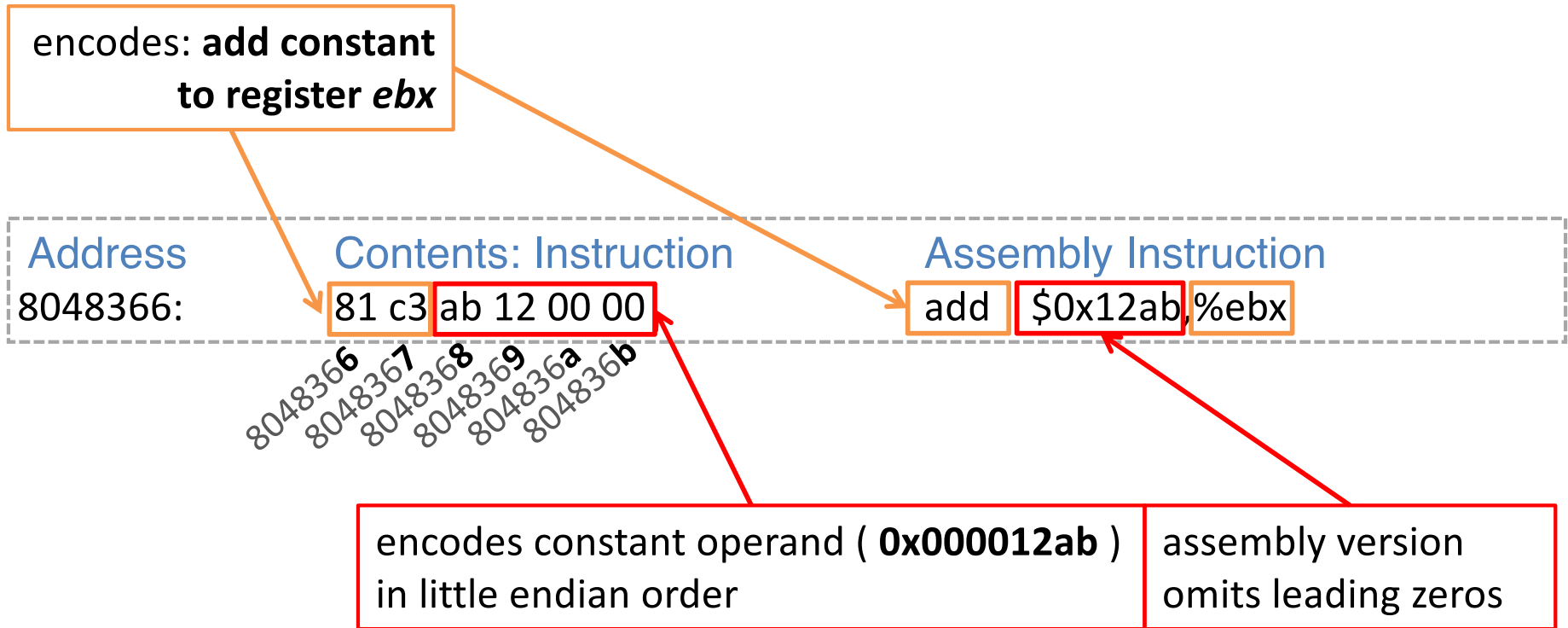## *Little Endian:* least significant byte first
- low order byte at low address, high order byte at high address
- used by **x86**, …

## *Big Endian:* most significant byte first
- high order byte at low address, low order byte at high address
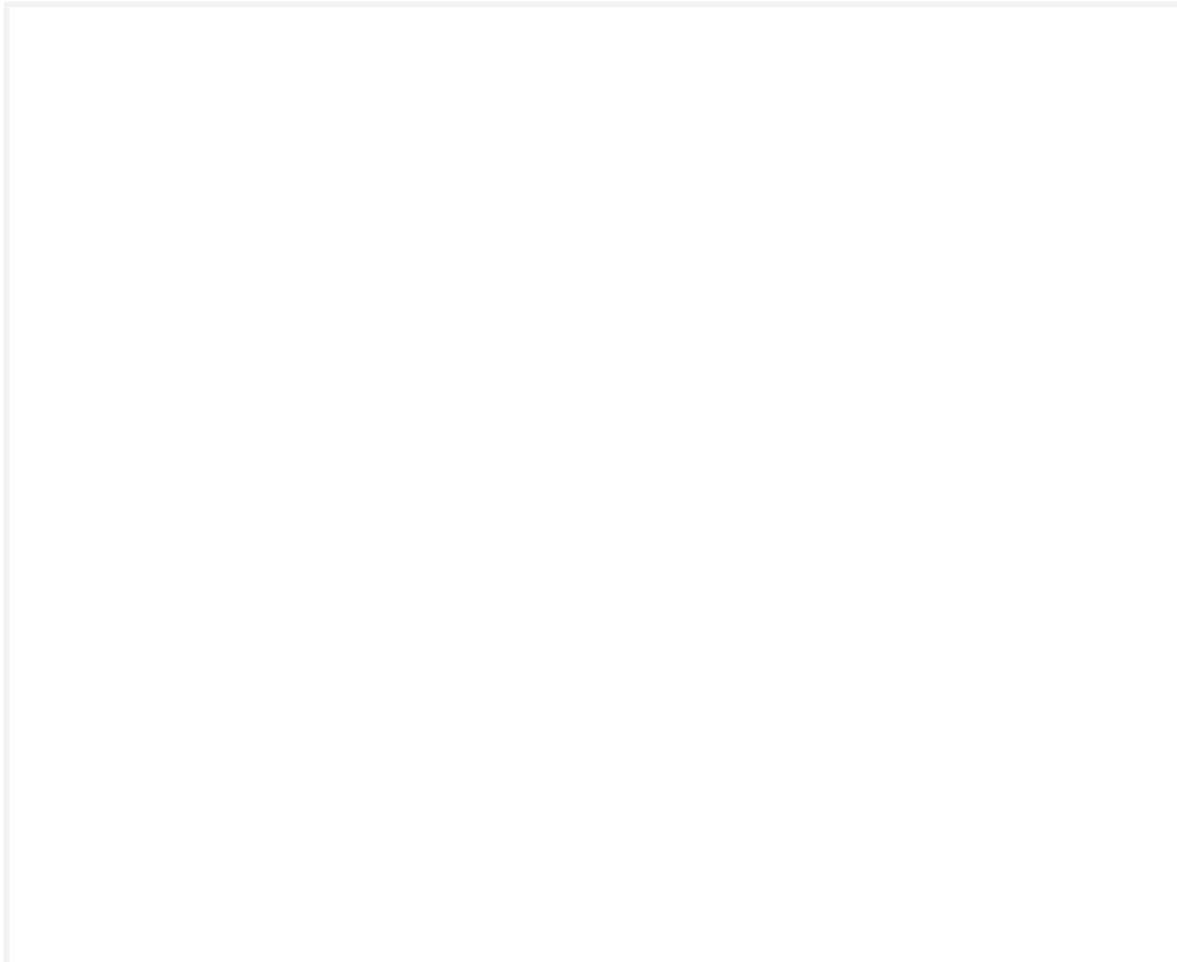- used by networks, SPARC, …

# Endianness in Machine Code

encodes: **add constant to register *ebx***

Address
8048366:

Contents: Instruction
81 c3 | ab 12 00 00

8048366
8048367
8048368
8048369
804836a
804836b

Assembly Instruction
add | $0x12ab,%ebx

encodes constant operand ( **0x000012ab** )
in little endian order

assembly version
omits leading zeros

# Data, Addresses, and Pointers

*address* = index of a cell in memory

*pointer* = address represented as data

| 0x03 | 0x02 | 0x01 | 0x00 | |
|------|------|------|------|------|
|      |      |      |      | 0x24 |
| 00   | 00   | 00   | F0   | 0x20 |
|      |      |      |      | 0x1C |
|      |      |      |      | 0x18 |
|      |      |      |      | 0x14 |
| 00   | 00   | 00   | 0C   | 0x10 |
|      |      |      |      | 0x0C |
| 00   | 00   | 00   | 20   | 0x08 |
|      |      |      |      | 0x04 |
| 00   | 00   | 00   | 08   | 0x00 |

memory drawn as 32-bit values, little endian order

# C: variables are memory locations (for now)

Compiler maps variable → memory location.

Declarations do not initialize!

```
int x; // x at 0x20
int y; // y at 0x0C

x = 0; // store 0 at 0x20

// store 0x3CD02700 at 0x0C
y = 0x3CD02700;

// load the contents at 0x0C,
// add 3, and store sum at 0x20
x = y + 3;
```
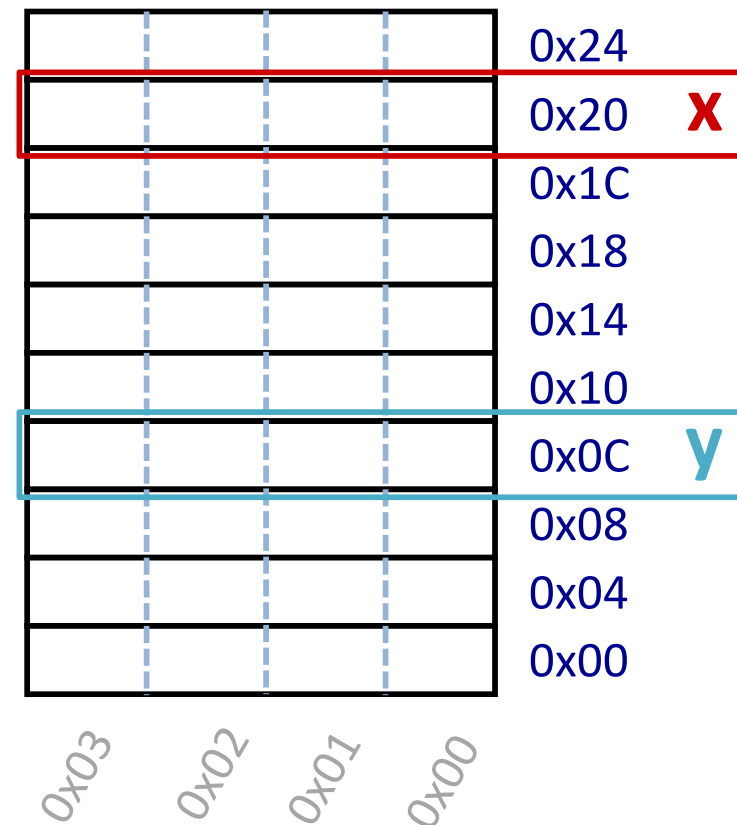
| | | | | 0x24 |
|---|---|---|---|---|
| | | | | 0x20  **x** |
| | | | | 0x1C |
| | | | | 0x18 |
| | | | | 0x14 |
| | | | | 0x10 |
| | | | | 0x0C  **y** |
| | | | | 0x08 |
| | | | | 0x04 |
| | | | | 0x00 |

0x03   0x02   0x01   0x00

# C: Address and Pointer Primitives

*address* = index of a cell/location in memory

*pointer* = address represented as data

**Expressions using addresses and pointers:**

&___     **address of** the memory location representing ___

*___     **contents at** the memory address given by ___

         a.k.a. "dereference ___"

**Pointer types:**

___*     address of a memory location holding a ___

# C: Address and Pointer Example

*& = address of*

*\* = contents at*

```
int* p;



int x = 5;
int y = 2;



p = &x;




y = 1 + *p;
```

18

# C: Address and Pointer Example

*& = address of*
*\* = contents at*

Declare a variable, `p`

```
int* p;
```

that will hold the address of a memory location holding an int

```
int x = 5;
int y = 2;
```

Declare two variables, `x` and `y`, that hold ints, and store 5 and 2 in them, respectively.

Get the address of the memory location

```
p = &x;
```

representing x

... and store it in p.     Now, "*p points to x.*"

Add 1 to the contents of memory at the address

```
y = 1 + *p;
```

stored in p

... and store it in the memory location representing y.

19

# C: Address and Pointer Example

*& = address of*
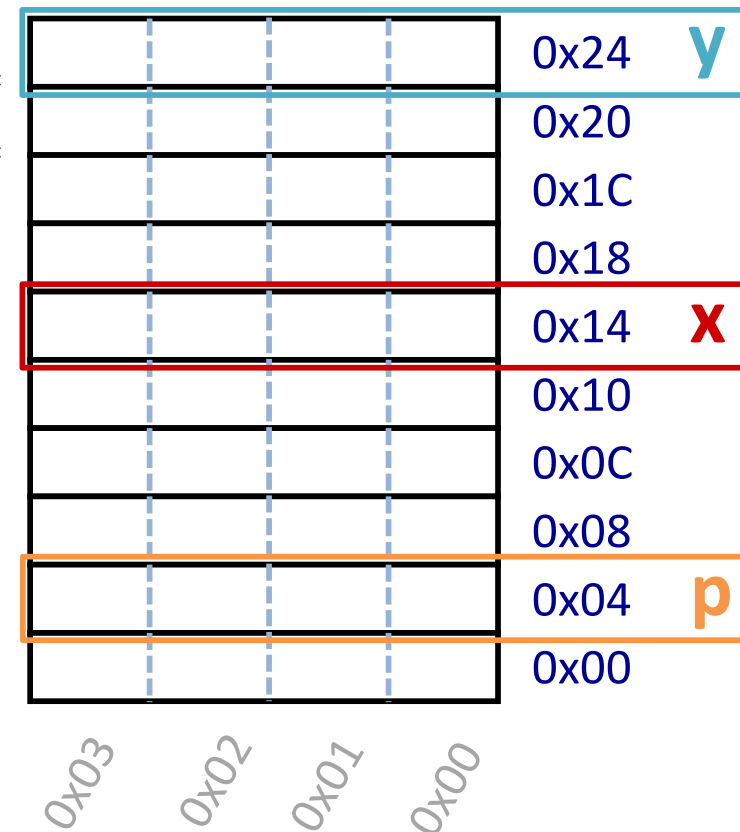*\* = contents at*

*C assignment:*

*Left-hand-side = right-hand-side;*

location

value

```
int* p;      // p: 0x04
int x = 5;   // x: 0x14, store 5 at 0x14
int y = 2;   // y: 0x24, store 2 at 0x24
p = &x;      // store 0x14 at 0x04
// load the contents at 0x04  (0x14)
// load the contents at 0x14  (0x5)
// add 1 and store sum at 0x24
y = 1 + *p;
// load the contents at 0x04  (0x14)
// store 0xF0 (240) at 0x14
*p = 240;
```

| | | | | |
|---|---|---|---|---|
| | | | | 0x24  y |
| | | | | 0x20 |
| | | | | 0x1C |
| | | | | 0x18 |
| | | | | 0x14  x |
| | | | | 0x10 |
| | | | | 0x0C |
| | | | | 0x08 |
| | | | | 0x04  p |
| | | | | 0x00 |

0x03   0x02   0x01   0x00

# C: Pointer Type Syntax

Spaces between base type, *, and variable name mostly do not matter.
**The following are <span style="color:red">equivalent</span>:**

```
int* ptr;
```
**I prefer this**

    I see: "The variable **ptr** holds an **address of an int** in memory."

```
int * ptr;
```

```
int *ptr;
```
**more common C style**

    I see: "Dereferencing the variable **ptr** will yield an **int**."

  Or   "The **memory location** where the variable **ptr** points holds an **int**."

    Caveat: do not declare multiple variables unless using the last form.
    `int* a, b;` means `int *a, b;` means `int* a; int b;`
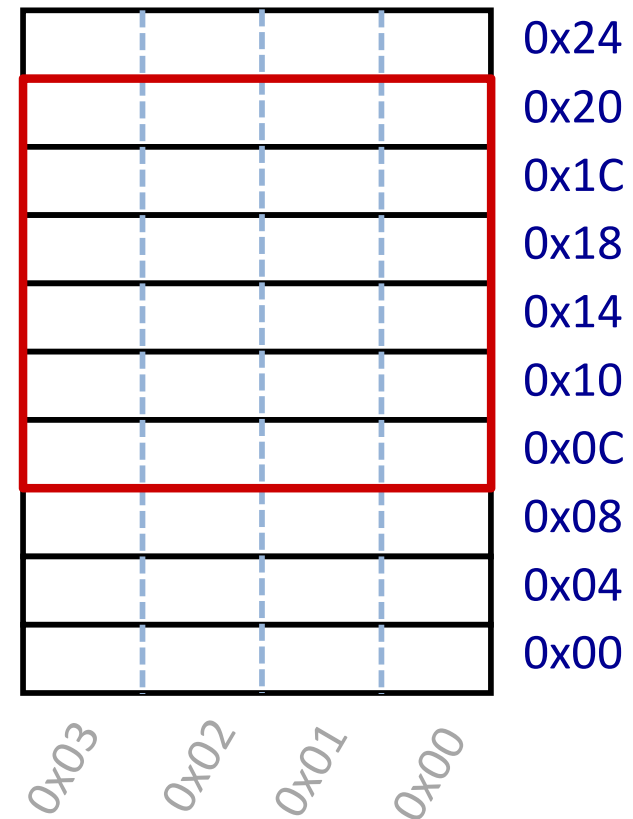
# C: Arrays

Declaration: `int a[6];`

element type

name

number of elements

Arrays are adjacent memory locations storing the same type of data.

**a** is a name for the array's base address, can be used as an *immutable* pointer.

| | 0x03 | 0x02 | 0x01 | 0x00 | |
|---|---|---|---|---|---|
| | | | | | 0x24 |
| | | | | | 0x20 |
| | | | | | 0x1C |
| | | | | | 0x18 |
| | | | | | 0x14 |
| | | | | | 0x10 |
| | | | | | 0x0C |
| | | | | | 0x08 |
| | | | | | 0x04 |
| | | | | | 0x00 |

# C: Arrays

Arrays are adjacent memory locations storing the same type of data.

**a** is a name for the array's base address, can be used as an *immutable* pointer.

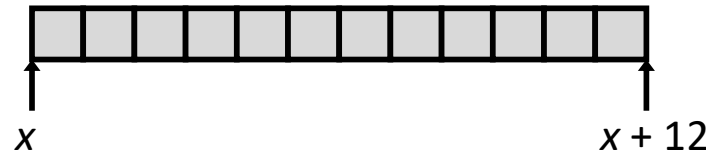Address of **a[i]** is base address **a** plus **i** times element size in bytes.

Declaration:
```
int a[6];
```

Indexing:
```
a[0] = 0xf0;
a[5] = a[0];
```

No bounds check:
```
a[6] = 0xBAD;
a[-1] = 0xBAD;
```

Pointers:
```
int* p;
```
equivalent {
```
p = a;
p = &a[0];
```
```
*p = 0xA;
```

equivalent {
```
p[1] = 0xB;
*(p + 1) = 0xB;
```
```
p = p + 2;
```

*array indexing = address arithmetic*
Both are scaled by the size of the type.

```
*p = a[1] + 1;
```

# C: Array Allocation

**Basic Principle**

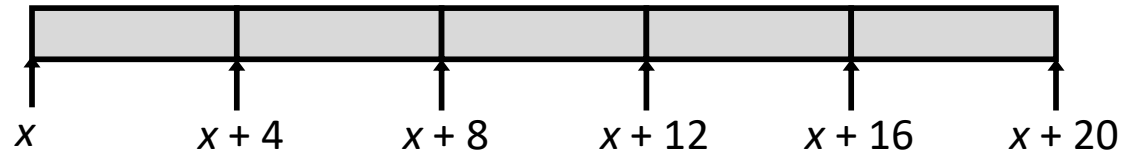$T \quad A[N];$

Array of length $N$ with elements of type $T$ and name $A$

*Contiguous* block of $N*sizeof(T)$ bytes of memory

`char string[12];`

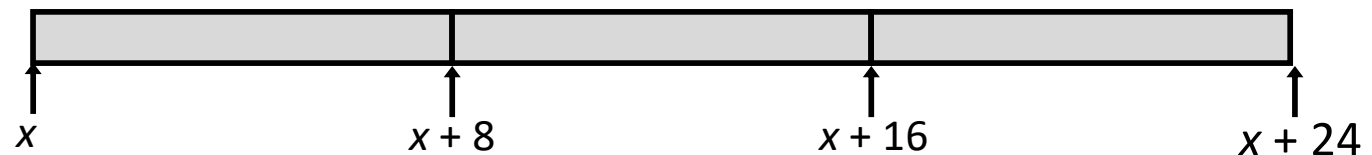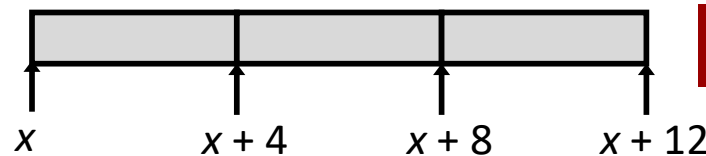Use *sizeof* to determine proper size in C.

$x$         $x + 12$

`int val[5];`

$x$    $x + 4$    $x + 8$    $x + 12$    $x + 16$    $x + 20$

`double a[3];`

$x$      $x + 8$      $x + 16$      $x + 24$

`char* p[3];`
(or `char *p[3];`)

IA32

$x$    $x + 4$    $x + 8$    $x + 12$
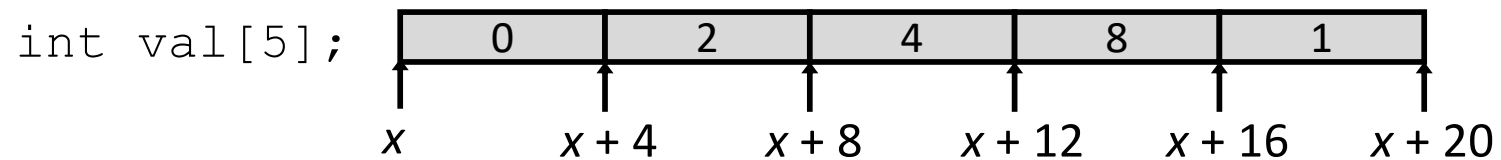
x86-64

$x$      $x + 8$      $x + 16$      $x + 24$

# C: Array Access

**Basic Principle**

$T$   $A[N]$;

Array of length $N$ with elements of type $T$ and name $A$

Identifier $A$ has type

int val[5];

| 0 | 2 | 4 | 8 | 1 |
|---|---|---|---|---|

$x$          $x + 4$          $x + 8$          $x + 12$          $x + 16$          $x + 20$

| Reference | Type | Value |
|-----------|------|-------|
| val[4] | int | |
| val | int * | |
| val+1 | int * | |
| &val[2] | int * | |
| val[5] | int | |
| *(val+1) | int | |
| val + i | int * | |

# C: **Null-terminated strings**

**ex**

C strings: arrays of ASCII characters ending with *null* character.

**Why?**

| 0x48 | 0x61 | 0x72 | 0x72 | 0x79 | 0x20 | 0x50 | 0x6F | 0x74 | 0x74 | 0x65 | 0x72 | 0x00 |
|------|------|------|------|------|------|------|------|------|------|------|------|------|
| 'H'  | 'a'  | 'r'  | 'r'  | 'y'  | ' '  | 'P'  | 'o'  | 't'  | 't'  | 'e'  | 'r'  | '\0' |

Does Endianness matter for strings?

```
int string_length(char str[]) {


}
```

# C: * *and* []

**C programmers often use * where you might expect []:**

   *e.g.,* char*:

   •    pointer to a char

   •    pointer to the first char in a string of unknown length

```
int strcmp(char* a, char* b);
int string_length(char* str) {
  // Try with pointer arithmetic, but no array indexing.




}
```
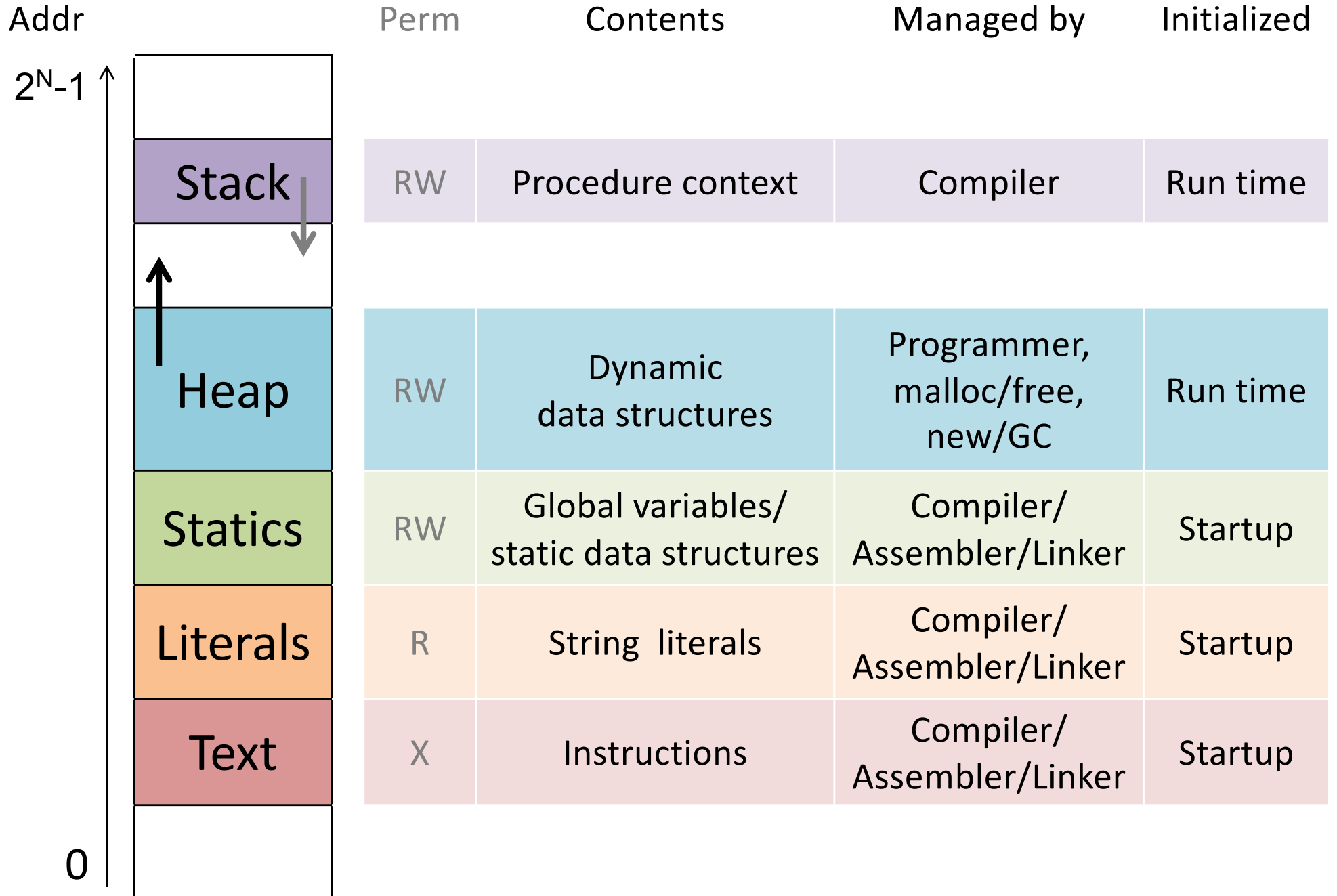
# NULL vs. '\0' vs. 0

**The null character (i.e., '\0' )** is used to signify the end of a string of characters.  It is one byte.  Equal to 0b00000000.

There is also **NULL.**  NULL is a pointer to an invalid memory address.  NULL is often used to indicate the end of an array of pointers.  Dereferencing a NULL pointer leads to crash!  In x86-64 machines, NULL is 8 bytes (same as the size of a pointer), equivalent to 0x00000000.
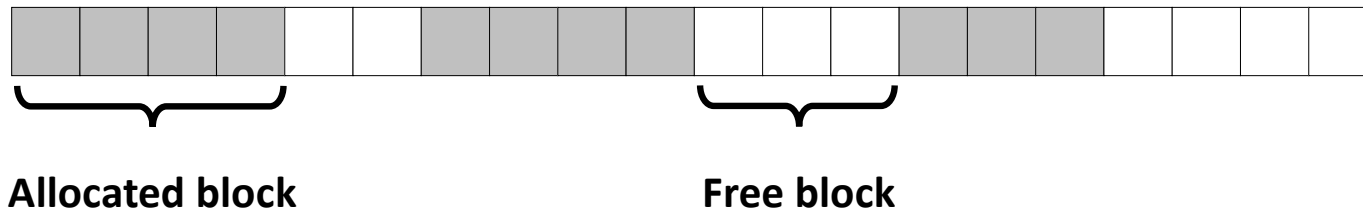
**0 is 0.**  It so happens that the null character and NULL are also literals of value 0.  NULL and the null character do not need to be equivalent to 0 but are done so by convention.

# Memory Layout

| Addr | | Perm | Contents | Managed by | Initialized |
|---|---|---|---|---|---|
| $2^N-1$ | | | | | |
| | Stack | RW | Procedure context | Compiler | Run time |
| | Heap | RW | Dynamic data structures | Programmer, malloc/free, new/GC | Run time |
| | Statics | RW | Global variables/ static data structures | Compiler/ Assembler/Linker | Startup |
| | Literals | R | String literals | Compiler/ Assembler/Linker | Startup |
| | Text | X | Instructions | Compiler/ Assembler/Linker | Startup |
| 0 | | | | | |

# C: Dynamic memory allocation in the heap

**Heap:**



**Allocated block**                    **Free block**

## Managed by memory allocator:

pointer to newly allocated block
of at least that size

number of contiguous bytes required

```
void* malloc(size_t size);
```

pointer to allocated block to free

```
void free(void* ptr);
```

# C: Dynamic array allocation

```
#define ZIP_LENGTH 5
int* zip = (int*)malloc(sizeof(int)*ZIP_LENGTH);
if (zip == NULL) {      // if error occurred
  perror("malloc");     // print error message
  exit(0);              // end the program
}

zip[0] = 0;
zip[1] = 2;
zip[2] = 4;
zip[3] = 8;
zip[4] = 1;


printf("zip is");
for (int i = 0; i < ZIP_LENGTH; i++) {
    printf(" %d", zip[i]);
}
printf("\n");

free(zip);
```
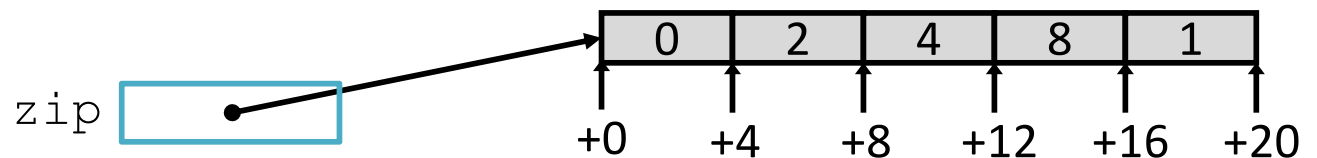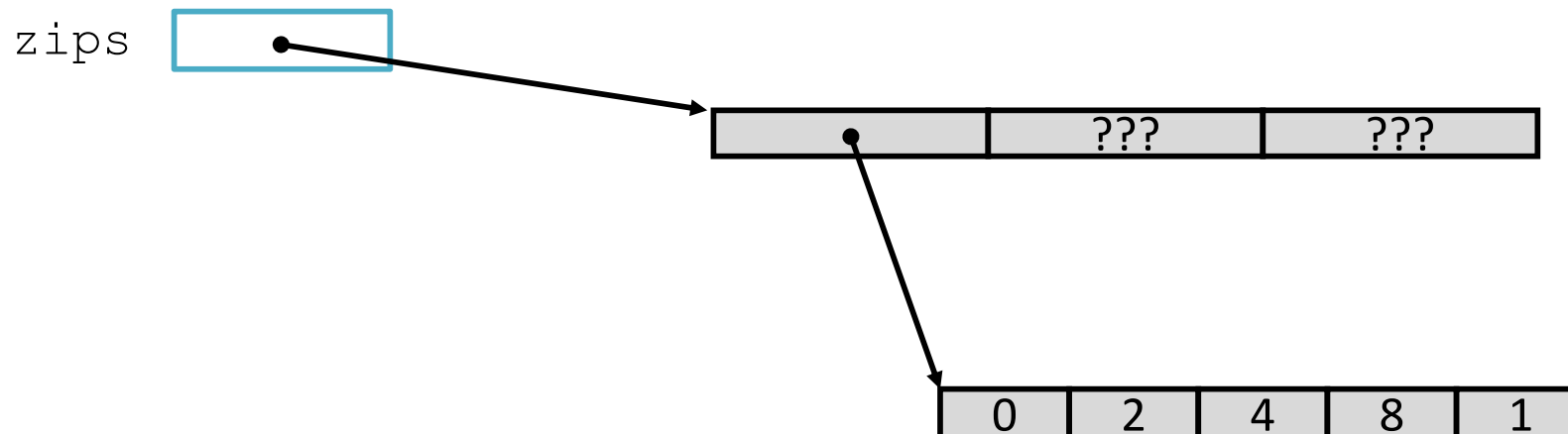
zip   0x7fedd2400dc0   0x7fff58bdd938

| | |
|---|---|
| 1 | 0x7fedd2400dd0 |
| 8 | 0x7fedd2400dcc |
| 4 | 0x7fedd2400dc8 |
| 2 | 0x7fedd2400dc4 |
| 0 | 0x7fedd2400dc0 |

zip

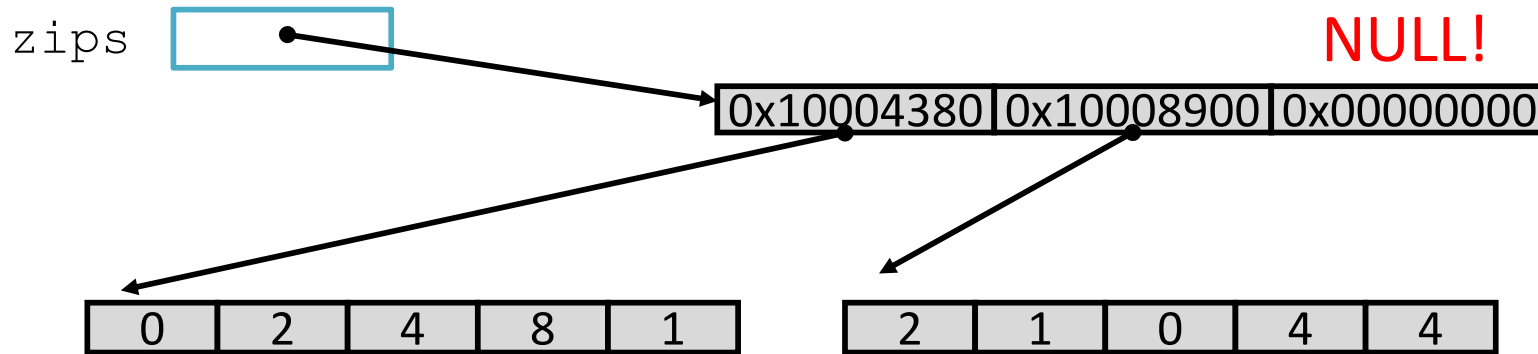| 0 | 2 | 4 | 8 | 1 |
|---|---|---|---|---|
| +0 | +4 | +8 | +12 | +16 | +20 |

# C: Arrays of pointers to arrays of …

```
int** zips = (int**)malloc(sizeof(int*)*3);
...
zips[0] = (int*)malloc(sizeof(int)*5);
...
int* zip0 = zips[0];
zip0[0] = 0;
zips[0][1] = 2;
zips[0][2] = 4;
zips[0][3] = 8;
zips[0][4] = 1;
```

zips

| | ??? | ??? |
|---|---|---|

| 0 | 2 | 4 | 8 | 1 |
|---|---|---|---|---|

# Array of Zip Codes

zips

| 0x10004380 | 0x10008900 | 0x00000000 |

NULL!

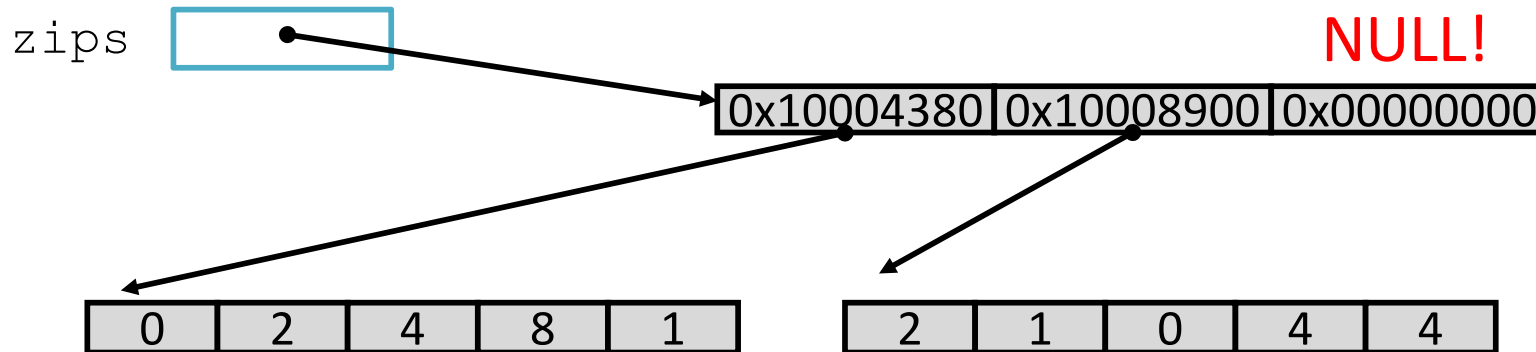| 0 | 2 | 4 | 8 | 1 |

| 2 | 1 | 0 | 4 | 4 |

Here we can use **NULL** to signify the end of an array of zip codes.

Why might it be important to end my array of pointers with NULL?  What if we wanted a function to print out all the zip codes in the array but we didn't know how long the array was?

What about ending zip codes with NULL?  Does that make sense?

# Zip Cycles

zips ⬚ → | 0x10004380 | 0x10008900 | 0x00000000 |

NULL!

| 0 | 2 | 4 | 8 | 1 |

| 2 | 1 | 0 | 4 | 4 |

```
// return a count of all zips the end with digit endNum

int zipCount(int* zips[], int endNum) {




}
```
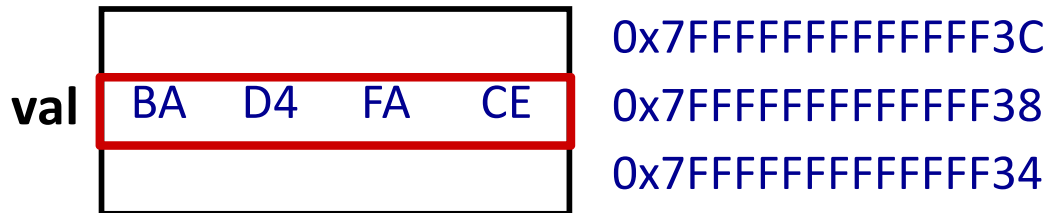
# C: scanf reads formatted input
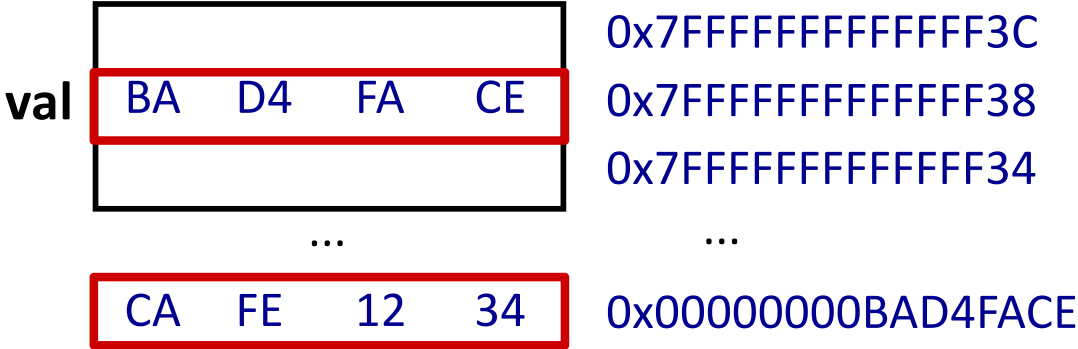
```
int val;

...

scanf("%d", &val);
```

Declared, but not initialized – holds anything.

Read one int from input.

Store it in memory at this address.

i.e., store it in memory at the address **where the contents of val is stored**: store into memory at 0xFFFFFF38.

| | |
|---|---|
| | 0x7FFFFFFFFFFFF3C |
| val  BA   D4   FA   CE | 0x7FFFFFFFFFFFF38 |
| | 0x7FFFFFFFFFFFF34 |

# C: classic bug using scanf

!!!

```
int val;
```
Declared, but not initialized – holds anything.

```
...

scanf("%d", val);
```

Read one int from input.

Store it in memory at this address.

i.e., store it in memory at the address **given by the contents of val**: store into memory at 0xBAD4FACE.

| | | | | |
|---|---|---|---|---|
| | | | | 0x7FFFFFFFFFFFFF3C |
| **val** | BA | D4 | FA | CE | 0x7FFFFFFFFFFFFF38 |
| | | | | 0x7FFFFFFFFFFFFF34 |
| | ... | | | ... |
| | CA | FE | 12 | 34 | 0x00000000BAD4FACE |

**Best case:** segmentation fault, or bus error, crash.

**Bad case:** silently corrupt data stored at address 0xBAD4FACE, and val still holds 0xBAD4FACE.
**Worst case:** arbitrary corruption

# C: memory error messages



http://xkcd.com/371/

**11: segmentation fault** ("segfault", SIGSEGV)

   accessing address outside legal area of memory

**10: bus error**

   accessing misaligned or other problematic address

More to come on debugging!

# C: Why?

## Why learn C?

- Think like actual computer (abstraction close to machine level) without dealing with machine code.
- Understand just how much Your Favorite Language provides.
- Understand just how much Your Favorite Language might cost.
- Classic.
- Still (more) widely used (than it should be).
- Pitfalls still fuel devastating reliability and security failures today.

## Why not use C?

- Probably not the right language for your next personal project.
- It "gets out of the programmer's way"
  even when the programmer is unwittingly running toward a cliff.
- Many advances in programming language design since then have produced languages that fix C's problems while keeping strengths.