# Executables & Arrays

CSE 351 Summer 2020

**Instructor:**

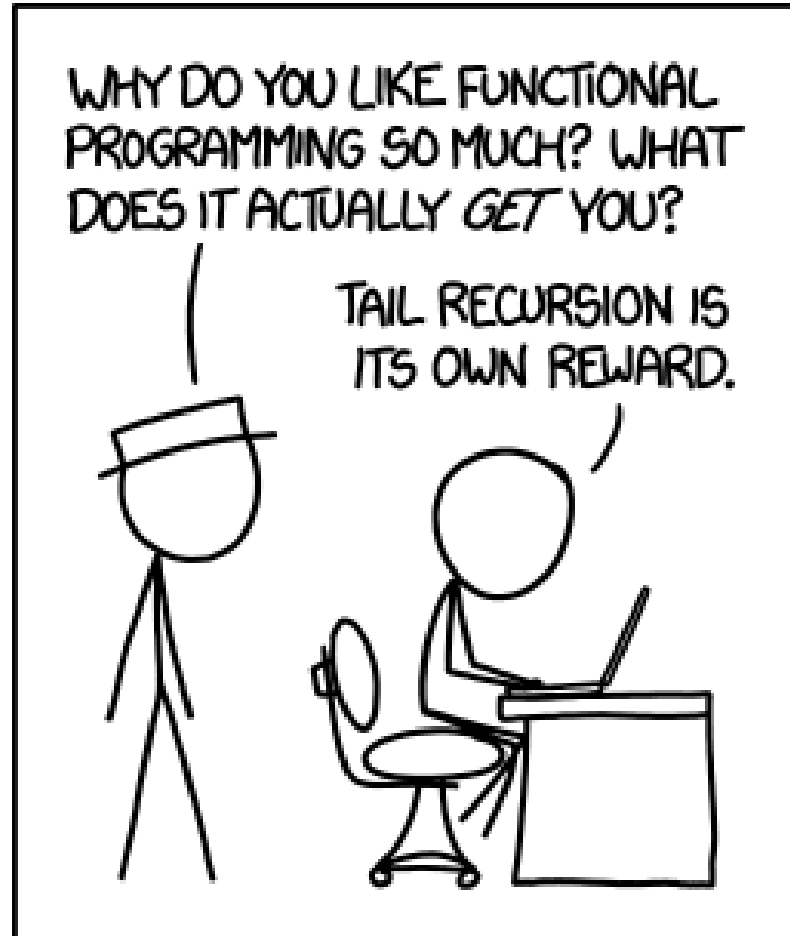Porter Jones

**Teaching Assistants:**

Amy Xu

Callum Walker

Sam Wolfson

Tim Mandzyuk



http://xkcd.com/1270/

# Administrivia

- Questions doc: https://tinyurl.com/CSE351-7-20

- hw12 due Wednesday (7/22) – 10:30am
- No hw due Friday!
- Lab 2 due Wednesday (7/22)
  - GDB Tutorial on Gradescope walks through first phase
  - Extra Credit portion – make sure you also submit to the Lab 2 Extra Credit assignment on Gradescope
- Thank you for the mid-quarter feedback!
  - Still sifting through it, will email with a summary soon
  - Can always provide anonymous feedback at https://feedback.cs.washington.edu

# Roadmap

C:

```
car *c = malloc(sizeof(car));
c->miles = 100;
c->gals = 17;
float mpg = get_mpg(c);
free(c);
```

Java:

```
Car c = new Car();
c.setMiles(100);
c.setGals(17);
float mpg =
        c.getMPG();
```

Memory & data
Integers & floats
x86 assembly
Procedures & stacks
Executables
Arrays & structs
Memory & caches
Processes
Virtual memory
Memory allocation
Java vs. C

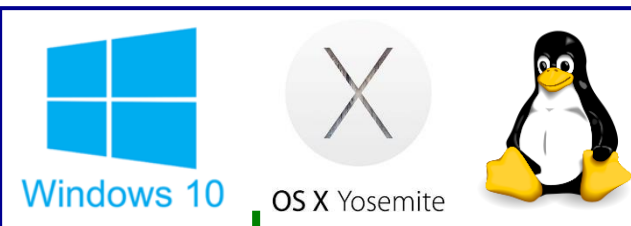Assembly
language:

```
get_mpg:
    pushq    %rbp
    movq     %rsp, %rbp
    ...
    popq     %rbp
    ret
```
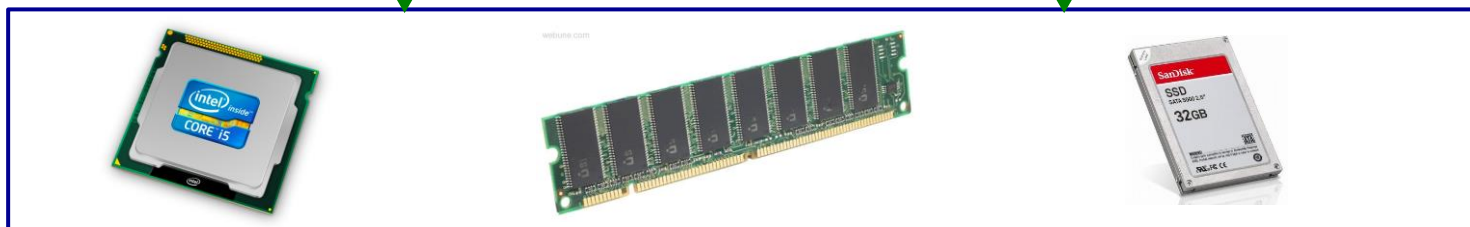
OS:

Machine
code:

```
0111010000011000
100011010000010000000010
1000100111000010
11000001111101000011111
```
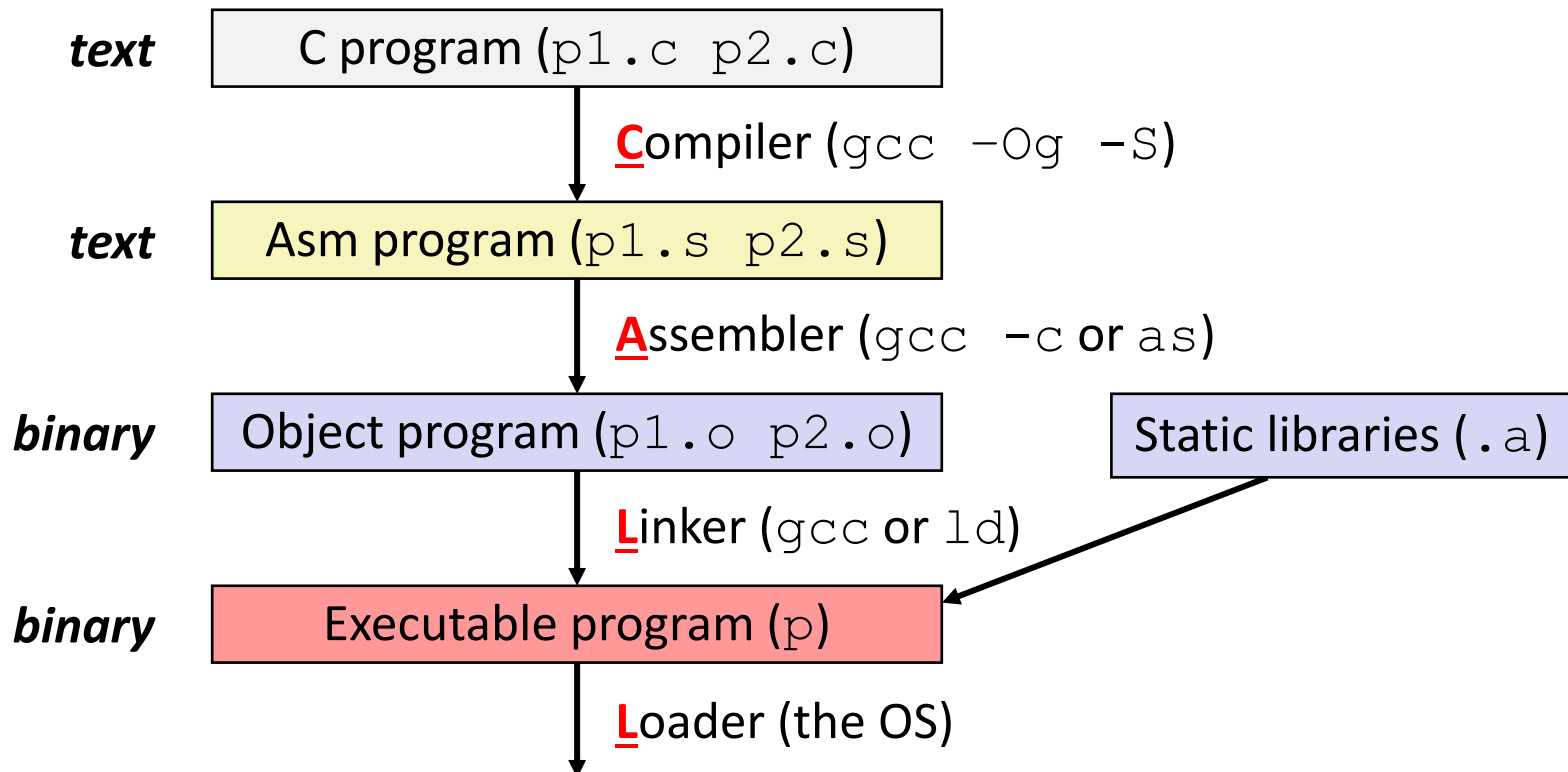
Windows 10   OS X Yosemite

Computer
system:

# Building an Executable from a C File

❖ Code in files `p1.c p2.c`

❖ Compile with command: `gcc -Og p1.c p2.c -o p`
  - Put resulting machine code in file `p`

❖ Run with command: `./p`

*text* | C program (`p1.c p2.c`)

↓ **C**ompiler (`gcc -Og -S`)

*text* | Asm program (`p1.s p2.s`)

↓ **A**ssembler (`gcc -c` or `as`)

*binary* | Object program (`p1.o p2.o`)          Static libraries (`.a`)

↓ **L**inker (`gcc` or `ld`)

*binary* | Executable program (`p`)

↓ **L**oader (the OS)

4

# Compiler

- ❖ **Input:** Higher-level language code (*e.g.* C, Java)
  - ▪ `foo.c`
- ❖ **Output:** Assembly language code (*e.g.* x86, ARM, MIPS)
  - ▪ `foo.s`

- ❖ First there's a preprocessor step to handle #directives
  - ▪ Macro substitution, plus other specialty directives
  - ▪ If curious/interested:  http://tigcc.ticalc.org/doc/cpp.html
- ❖ Super complex, whole courses devoted to these!
- ❖ Compiler optimizations
  - ▪ "Level" of optimization specified by capital 'O' flag (*e.g.* `-Og`, `-O3`)
  - ▪ Options:  https://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html

# Compiling Into Assembly

❖ C Code (`sum.c`)

```
void sumstore(long x, long y, long *dest) {
    long t = x + y;
    *dest = t;
}
```

❖ x86-64 assembly (`gcc –Og –S sum.c`)

```
sumstore(long, long, long*):
  addq     %rdi, %rsi
  movq     %rsi, (%rdx)
  ret
```

Warning:  You may get different results with other versions of `gcc` and different compiler settings

# Assembler

❖ **Input:** Assembly language code (*e.g.* x86, ARM, MIPS)
  - `foo.s`

❖ **Output:** Object files (*e.g.* ELF, COFF)
  - `foo.o`
  - Contains *object code* and *information tables*

❖ Reads and uses *assembly directives*
  - *e.g.* `.text, .data, .quad`
  - x86: https://docs.oracle.com/cd/E26502_01/html/E28388/eoiyg.html

❖ Produces "machine language"
  - Does its best, but object file is *not* a completed binary

❖ <u>Example</u>: `gcc -c foo.s`

# Producing Machine Language

- ❖ **Simple cases:** arithmetic and logical operations, shifts, etc.
    - All necessary information is contained in the instruction itself

- ❖ What about the following?
    - Conditional jump
    - Accessing static data (*e.g.* global var or jump table)
    - `call`

- ❖ Addresses and labels are problematic because the final executable hasn't been constructed yet!
    - So how do we deal with these in the meantime?

# Object File Information Tables

- ❖ **Symbol Table** holds list of "items" that may be used by other files
  - ▪ *Non-local labels* – function names for `call`
  - ▪ *Static Data* – variables & literals that might be accessed across files

- ❖ **Relocation Table** holds list of "items" that this file needs the address of later (currently undetermined)
  - ▪ Any *label* or piece of *static data* referenced in an instruction in this file
    - • Both internal and external

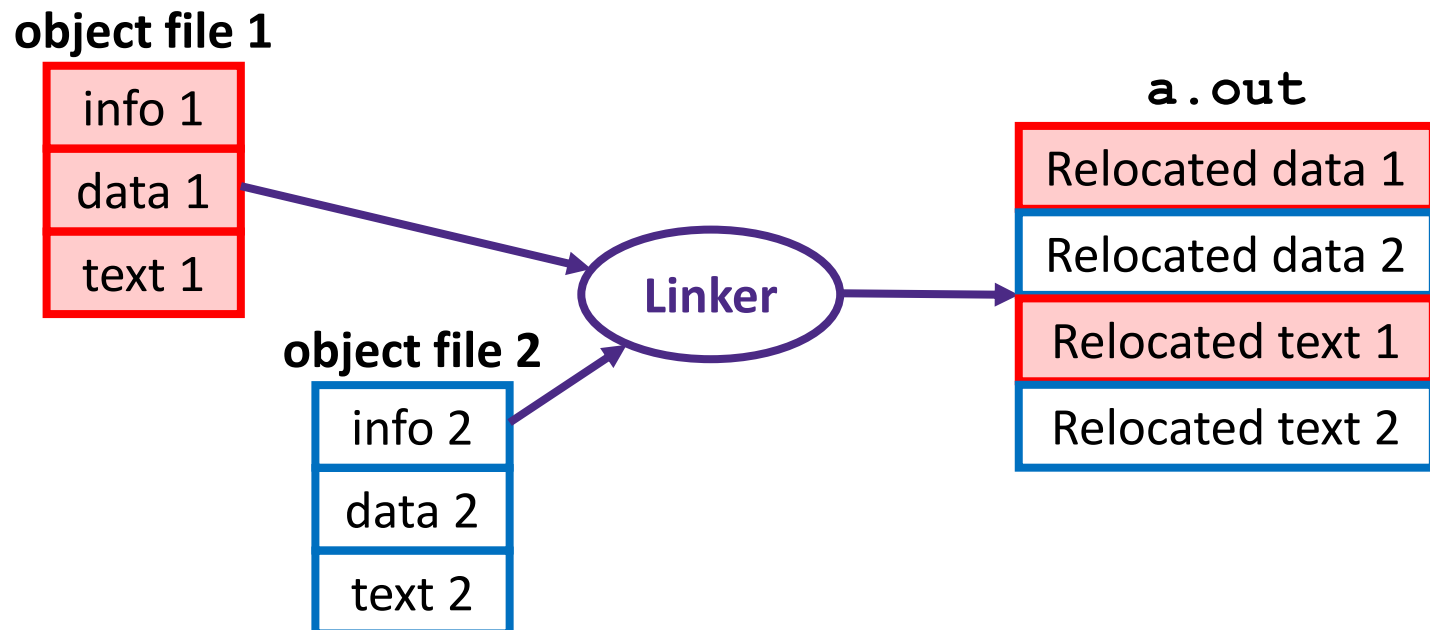- ❖ Each file has its own symbol and relocation tables

# Object File Format

1) <u>object file header</u>: size and position of the other pieces of the object file

2) <u>text segment</u>: the machine code

3) <u>data segment</u>: data in the source file (binary)

4) <u>relocation table</u>: identifies lines of code that need to be "handled"

5) <u>symbol table</u>: list of this file's labels and data that can be referenced

6) <u>debugging information</u>

❖ More info: ELF format
  ▪ http://www.skyfree.org/linux/references/ELF_Format.pdf

# Linker

- ❖ **Input:**  Object files (e.g. ELF, COFF)
  - ▪ `foo.o`
- ❖ **Output:**  executable binary program
  - ▪ `a.out`

- ❖ Combines several object files into a single executable (*linking*)
- ❖ Enables separate compilation/assembling of files
  - ▪ Changes to one file do not require recompiling of whole program

# Linking

1) Take text segment from each `.o` file and put them together
2) Take data segment from each `.o` file, put them together, and concatenate this onto end of text segments
3) Resolve References
   ▪ Go through Relocation Table; handle each entry

**object file 1**

| info 1 |
| --- |
| data 1 |
| text 1 |

**object file 2**

| info 2 |
| --- |
| data 2 |
| text 2 |

**Linker**

**`a.out`**

| Relocated data 1 |
| --- |
| Relocated data 2 |
| Relocated text 1 |
| Relocated text 2 |

# Disassembling Object Code

❖ Disassembled:

```
0000000000400536 <sumstore>:
  400536:   48 01 fe        add     %rdi,%rsi
  400539:   48 89 32        mov     %rsi,(%rdx)
  40053c:   c3              retq
```

❖ **Disassembler** (`objdump -d sum`)
  - Useful tool for examining object code (`man 1 objdump`)
  - Analyzes bit pattern of series of instructions
  - Produces approximate rendition of assembly code
  - Can run on either `a.out` (complete executable) or `.o` file

# What Can be Disassembled?

```
% objdump -d WINWORD.EXE

WINWORD.EXE:     file format pei-i386

No symbols in "WINWORD.EXE".
Disassembly of section .text:

30001000 <.text>:
30001000:
30001001:
30001003:
30001005:
3000100a:
```

**Reverse engineering forbidden by
Microsoft End User License Agreement**

- ❖ Anything that can be interpreted as executable code
- ❖ Disassembler examines bytes and attempts to reconstruct assembly source

# Loader

❖ **Input:**  executable binary program, command-line arguments
  ▪ `./a.out arg1 arg2`
❖ **Output:**  <program is run>

❖ Loader duties primarily handled by OS/kernel
  ▪ More about this when we learn about processes
❖ Memory sections (Instructions, Static Data, Stack) are set up
❖ Registers are initialized

# Roadmap

C:

```
car *c = malloc(sizeof(car));
c->miles = 100;
c->gals = 17;
float mpg = get_mpg(c);
free(c);
```

Java:

```
Car c = new Car();
c.setMiles(100);
c.setGals(17);
float mpg =
        c.getMPG();
```

Memory & data
Integers & floats
x86 assembly
Procedures & stacks
Executables
**Arrays & structs**
Memory & caches
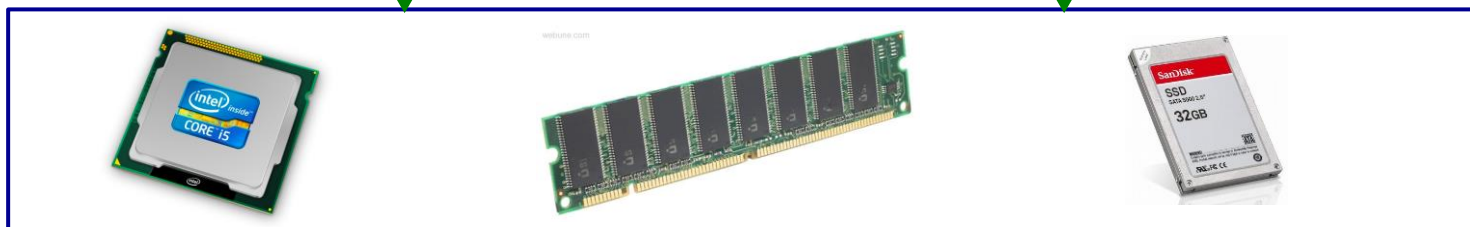Processes
Virtual memory
Memory allocation
Java vs. C

Assembly
language:

```
get_mpg:
    pushq    %rbp
    movq     %rsp, %rbp
    ...
    popq     %rbp
    ret
```
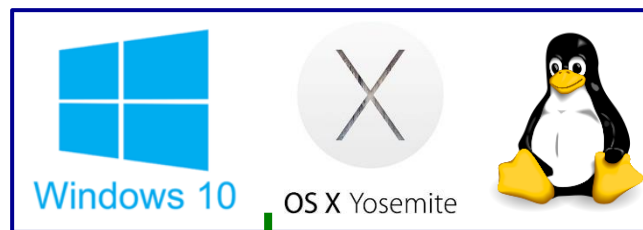
OS:

Machine
code:

```
0111010000011000
100011010000010000000010
1000100111000010
110000011111101000011111
```

Windows 10    OS X Yosemite

Computer
system:

intel inside CORE i5    SanDisk SSD 32GB

# Data Structures in Assembly

❖ **Arrays**
  ▪ **One-dimensional**
  ▪ Multidimensional (nested)
  ▪ Multilevel
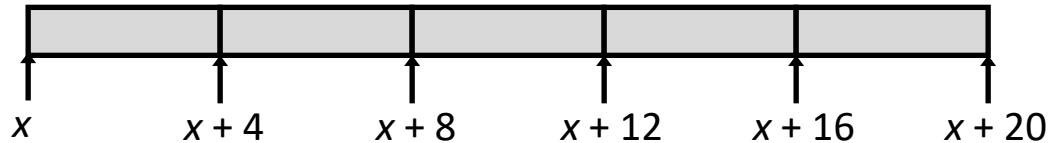
❖ Structs
  ▪ Alignment

❖ ~~Unions~~

17

# Review: Array Allocation

❖ Basic Principle
- **T** A[N];  →  array of data type **T** and length N
- *Contiguously* allocated region of N*sizeof(**T**) bytes
- Identifier A returns address of array (type **T\***)

# Review: Array Access

❖ **Basic Principle**
  - **T** A[N];    →    array of data type **T** and length N
  - Identifier A returns address of array (type **T\***)

```
int x[5];
```
|  3  |  7  |  1  |  9  |  5  |
| a | a+4 | a+8 | a+12 | a+16 | a+20 |

❖ <u>Reference</u>          <u>Type</u>          <u>Value</u>

| Reference | Type | Value |
| --- | --- | --- |
| x[4] | **int** | 5 |
| x | **int\*** | a |
| x+1 | **int\*** | a + 4 |
| &x[2] | **int\*** | a + 8 |
| x[5] | **int** | ??  (whatever's in memory at addr x+20) |
| *(x+1) | **int** | 7 |
| x+i | **int\*** | a + 4*i |

19

# Array Example

```
// arrays of ZIP code digits
int cmu[5] = { 1, 5, 2, 1, 3 };
int  uw[5] = { 9, 8, 1, 9, 5 };
int ucb[5] = { 9, 4, 7, 2, 0 };
```

brace-enclosed
list initialization

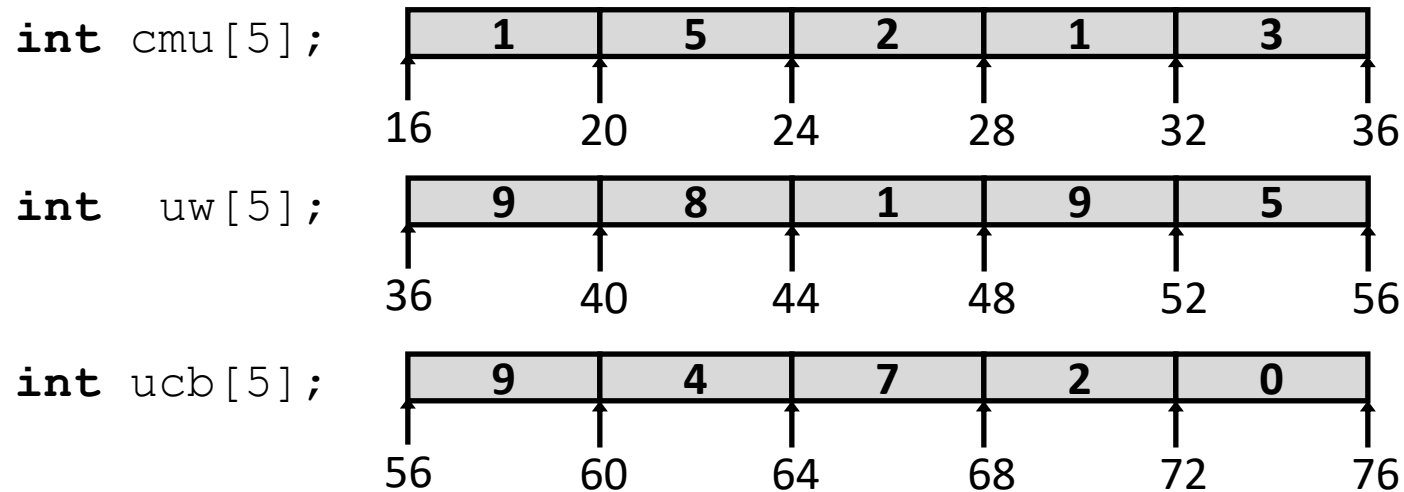# Array Example

```
// arrays of ZIP code digits
int cmu[5] = { 1, 5, 2, 1, 3 };
int  uw[5] = { 9, 8, 1, 9, 5 };
int ucb[5] = { 9, 4, 7, 2, 0 };
```

`int cmu[5];`

| 1 | 5 | 2 | 1 | 3 |
|---|---|---|---|---|

16    20    24    28    32    36

`int  uw[5];`

| 9 | 8 | 1 | 9 | 5 |
|---|---|---|---|---|

36    40    44    48    52    56

`int ucb[5];`

| 9 | 4 | 7 | 2 | 0 |
|---|---|---|---|---|

56    60    64    68    72    76

❖ Example arrays happened to be allocated in successive 20 byte blocks
  ▪ Not guaranteed to happen in general

# Array Accessing Example

```
int   uw[5];
```
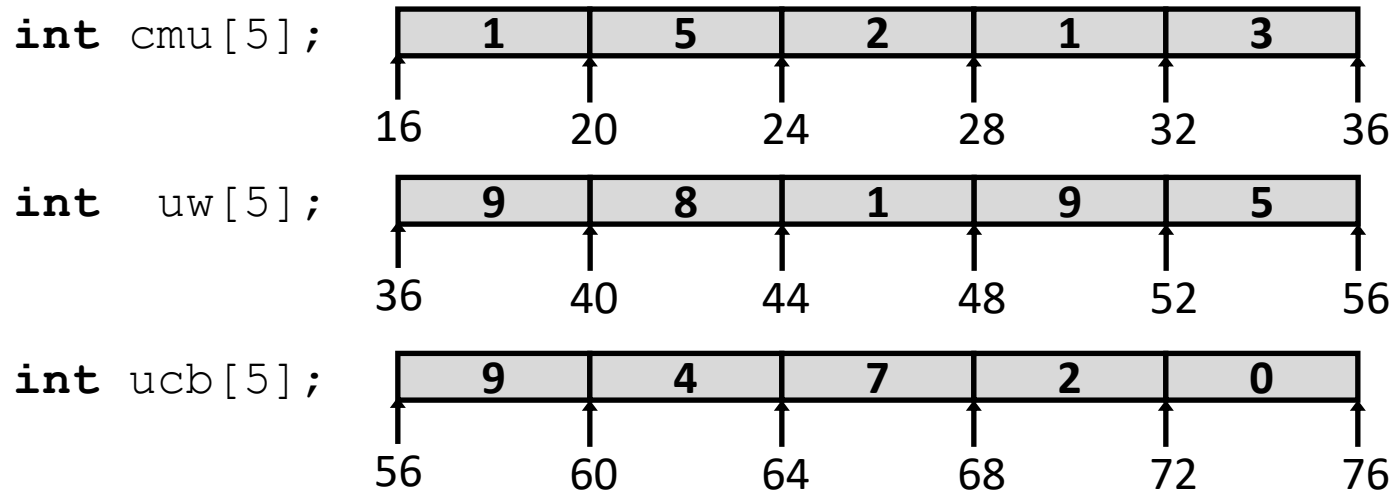
| 9 | 8 | 1 | 9 | 5 |
|---|---|---|---|---|

36    40    44    48    52    56

```
// return specified digit of ZIP code
int get_digit(int z[5], int digit) {
  return z[digit];
}
```

```
get_digit:
  movl (%rdi,%rsi,4), %eax   # z[digit]
```

- Register `%rdi` contains starting address of array

- Register `%rsi` contains array index

- Desired digit at `%rdi+4*%rsi`, so use memory reference `(%rdi,%rsi,4)`

# Referencing Examples

`int cmu[5];`

| 1 | 5 | 2 | 1 | 3 |
|---|---|---|---|---|

16      20      24      28      32      36

`int uw[5];`

| 9 | 8 | 1 | 9 | 5 |
|---|---|---|---|---|

36      40      44      48      52      56

`int ucb[5];`

| 9 | 4 | 7 | 2 | 0 |
|---|---|---|---|---|

56      60      64      68      72      76

| Reference | Address | Value | Guaranteed? |
|-----------|---------|-------|-------------|
| `uw[3]`   |         |       |             |
| `uw[6]`   |         |       |             |
| `uw[-1]`  |         |       |             |
| `cmu[15]` |         |       |             |

- ❖ No bounds checking
- ❖ Example arrays happened to be allocated in successive 20 byte blocks
  - ▪ Not guaranteed to happen in general

23

# C Details:  Arrays and Pointers

❖ **Arrays are (almost) identical to pointers**
  ▪ `char *string` **and** `char string[]` **are nearly identical declarations**
  ▪ Differ in subtle ways:  initialization, `sizeof()`, etc.

❖ **An array name is an expression (not a variable) that returns the address of the array**
  ▪ It *looks* like a pointer to the first (0[th]) element
    • `*ar` same as `ar[0]`, `*(ar+2)` same as `ar[2]`
  ▪ An array name is read-only (no assignment) because it is a *label*
    • Cannot use "`ar = <anything>`"

# C Details:  Arrays and Functions

❖ Declared arrays only allocated while the scope is valid:

```
char* foo() {
    char string[32]; ...;
    return string;
}
```

# BAD!

❖ An array is passed to a function as a pointer:

- Array size gets lost!

*Really* `int *ar`

```
int foo(int ar[], unsigned int size) {
    ... ar[size-1] ...
}
```

Must explicitly
pass the size!

UNIVERSITY *of* WASHINGTON

# Data Structures in Assembly

❖ **Arrays**
  ▪ One-dimensional
  ▪ **Multidimensional (nested)**
  ▪ Multilevel

❖ Structs
  ▪ Alignment

❖ Unions

# Nested Array Example

```
int sea[4][5] =
  {{ 9, 8, 1, 9, 5 },
   { 9, 8, 1, 0, 5 },
   { 9, 8, 1, 0, 3 },
   { 9, 8, 1, 1, 5 }};
```
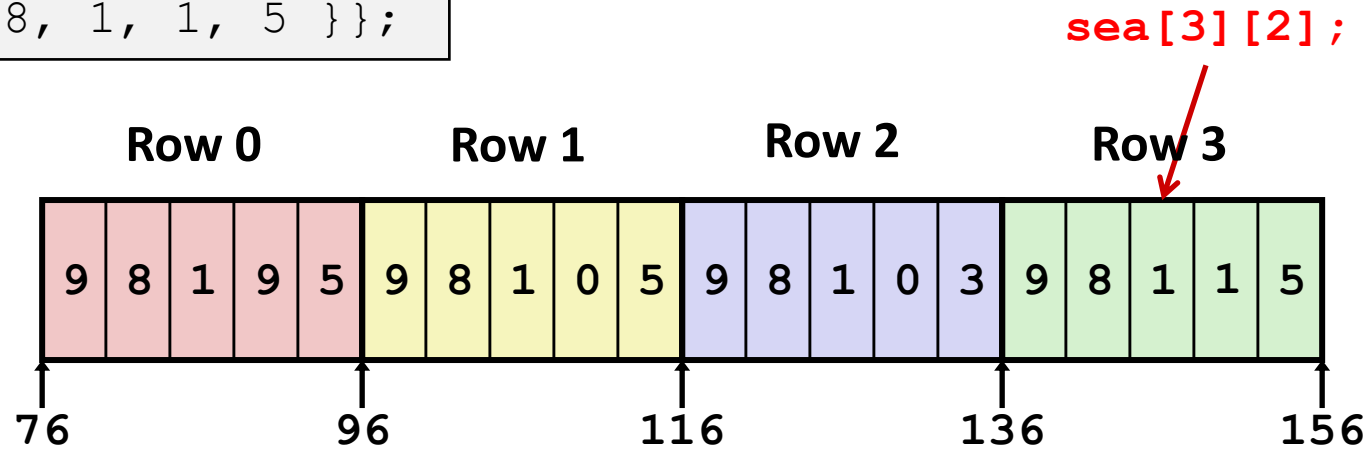
Remember, `T A[N]` is an array with elements of type `T`, with length `N`

❖ What is the layout in memory?

# Nested Array Example

```
int sea[4][5] =
  {{ 9, 8, 1, 9, 5 },
   { 9, 8, 1, 0, 5 },
   { 9, 8, 1, 0, 3 },
   { 9, 8, 1, 1, 5 }};
```

Remember, `T A[N]` is an array with elements of type `T`, with length `N`

`sea[3][2];`

| Row 0 | Row 1 | Row 2 | Row 3 |
|---|---|---|---|

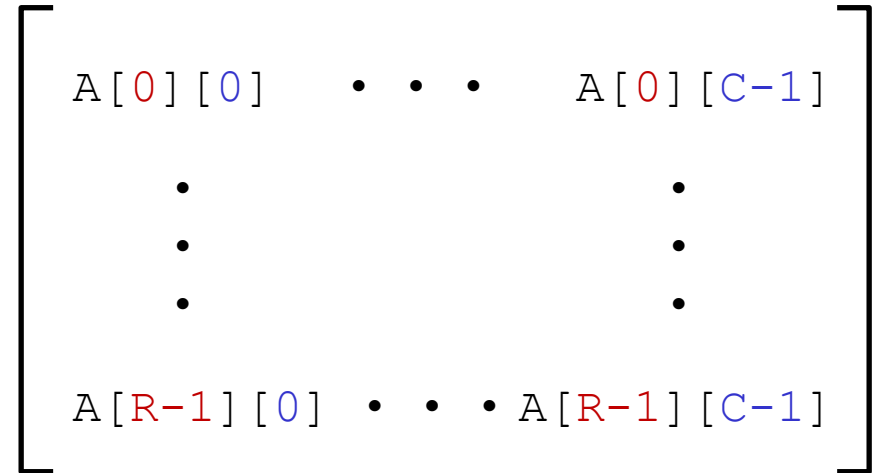| 9 | 8 | 1 | 9 | 5 | 9 | 8 | 1 | 0 | 5 | 9 | 8 | 1 | 0 | 3 | 9 | 8 | 1 | 1 | 5 |

76    96    116    136    156

- ❖ "Row-major" ordering of all elements
- ❖ Elements in the same row are contiguous
- ❖ Guaranteed (in C)

# Two-Dimensional (Nested) Arrays

❖ Declaration: `T A[R][C];`

- 2D array of data type `T`
- `R` rows, `C` columns
- Each element requires
  **sizeof(T)** bytes

❖ Array size?

$$
\begin{bmatrix}
\texttt{A[0][0]} & \cdots & \texttt{A[0][C-1]} \\
\vdots & & \vdots \\
\vdots & & \vdots \\
\vdots & & \vdots \\
\texttt{A[R-1][0]} & \cdots & \texttt{A[R-1][C-1]}
\end{bmatrix}
$$
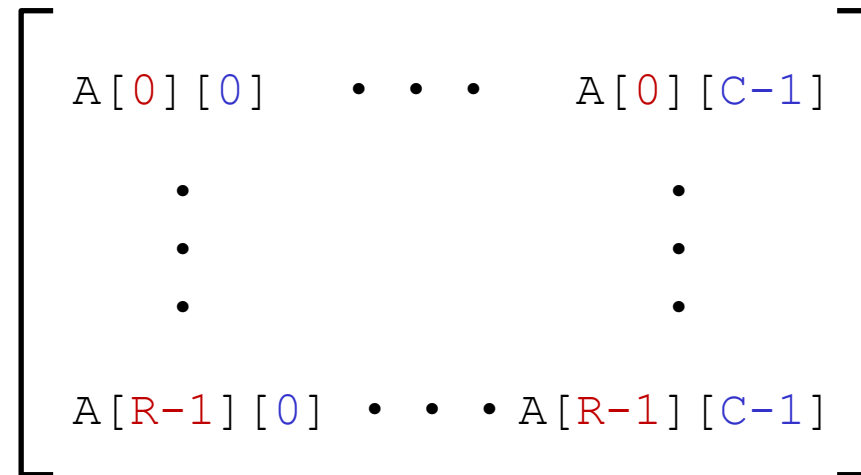
# Two-Dimensional (Nested) Arrays

❖ Declaration: `T A[R][C];`

- 2D array of data type `T`
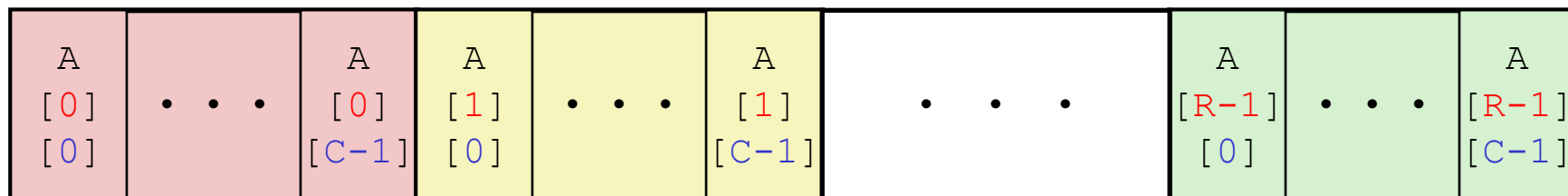- `R` rows, `C` columns
- Each element requires **`sizeof(T)`** bytes

$$\begin{bmatrix} A[0][0] & \cdots & A[0][C-1] \\ \vdots & & \vdots \\ \vdots & & \vdots \\ \vdots & & \vdots \\ A[R-1][0] & \cdots & A[R-1][C-1] \end{bmatrix}$$

❖ Array size:

- `R*C*`**`sizeof(T)`** bytes

❖ Arrangement: **row-major** ordering

`int A[R][C];`

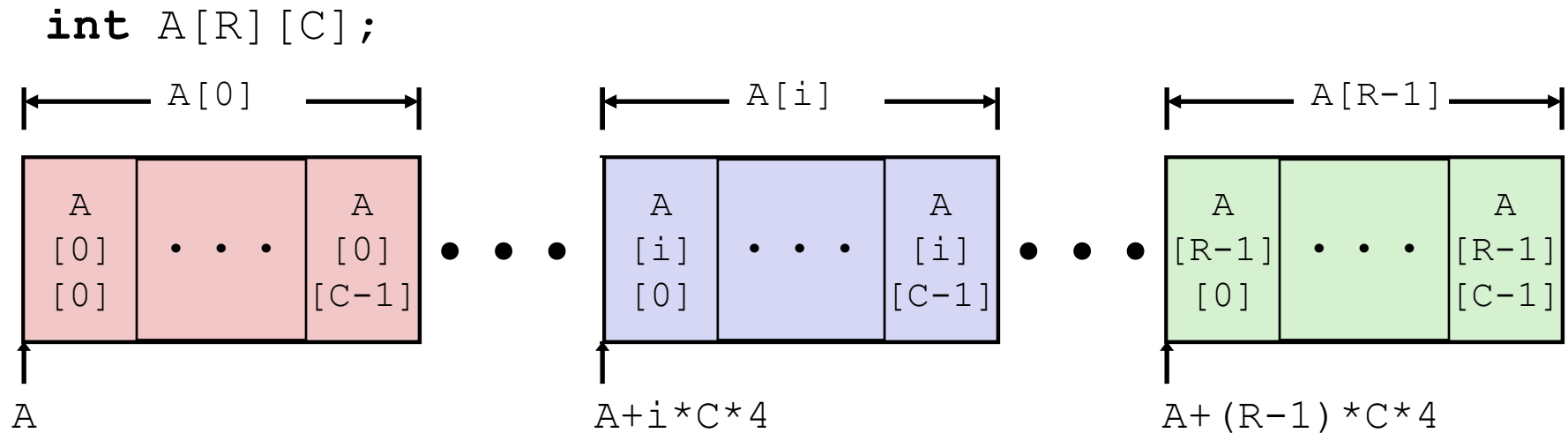| A [0] [0] | · · · | A [0] [C-1] | A [1] [0] | · · · | A [1] [C-1] | · · · | A [R-1] [0] | · · · | A [R-1] [C-1] |
|---|---|---|---|---|---|---|---|---|---|

←——————————————— `4*R*C` bytes ———————————————→

# **Nested Array <u>Row Access</u>**

❖ Row vectors

  ▪ Given **T** `A[R][C]`,

    • `A[i]` is an array of `C` elements ("row `i`")

    • `A` is address of array

    • Starting address of row `i` = `A + i*(C * sizeof(T))`

`int A[R][C];`



|  A[0]  |  |  | | | A[i] | | | | | A[R-1] | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A[0][0] | • • • | A[0][C-1] | ● ● ● | A[i][0] | • • • | A[i][C-1] | ● ● ● | A[R-1][0] | • • • | A[R-1][C-1] | | |

A                                    A+i*C*4                              A+(R-1)*C*4

# Nested Array <u>Row Access</u> Code

```
int* get_sea_zip(int index)
{
  return sea[index];
}
```

```
int sea[4][5] =
  {{ 9, 8, 1, 9, 5 },
   { 9, 8, 1, 0, 5 },
   { 9, 8, 1, 0, 3 },
   { 9, 8, 1, 1, 5 }};
```

```
get_sea_zip(int):
    movslq  %edi, %rdi
    leaq    (%rdi,%rdi,4), %rax
    leaq    sea(,%rax,4), %rax
    ret

sea:
    .long   9
    .long   8
    .long   1
    .long   9
    .long   5
    .long   9
    .long   8
...
```

# Nested Array <u>Row Access</u> Code

```
int* get_sea_zip(int index)
{
  return sea[index];
}
```

```
int sea[4][5] =
   {{ 9, 8, 1, 9, 5 },
    { 9, 8, 1, 0, 5 },
    { 9, 8, 1, 0, 3 },
    { 9, 8, 1, 1, 5 }};
```

- What data type is `sea[index]`?

- What is its value?

```
# %rdi = index
leaq (%rdi,%rdi,4),%rax          Translation?
leaq sea(,%rax,4),%rax
```

# Nested Array <u>Row Access</u> Code

```c
int* get_sea_zip(int index)
{
  return sea[index];
}
```

```c
int sea[4][5] =
  {{ 9, 8, 1, 9, 5 },
   { 9, 8, 1, 0, 5 },
   { 9, 8, 1, 0, 3 },
   { 9, 8, 1, 1, 5 }};
```

```
# %rdi = index
leaq (%rdi,%rdi,4),%rax   # 5 * index
leaq sea(,%rax,4),%rax    # sea + (20 * index)
```

❖ Row Vector
  ▪ `sea[index]` is array of 5 `int`s
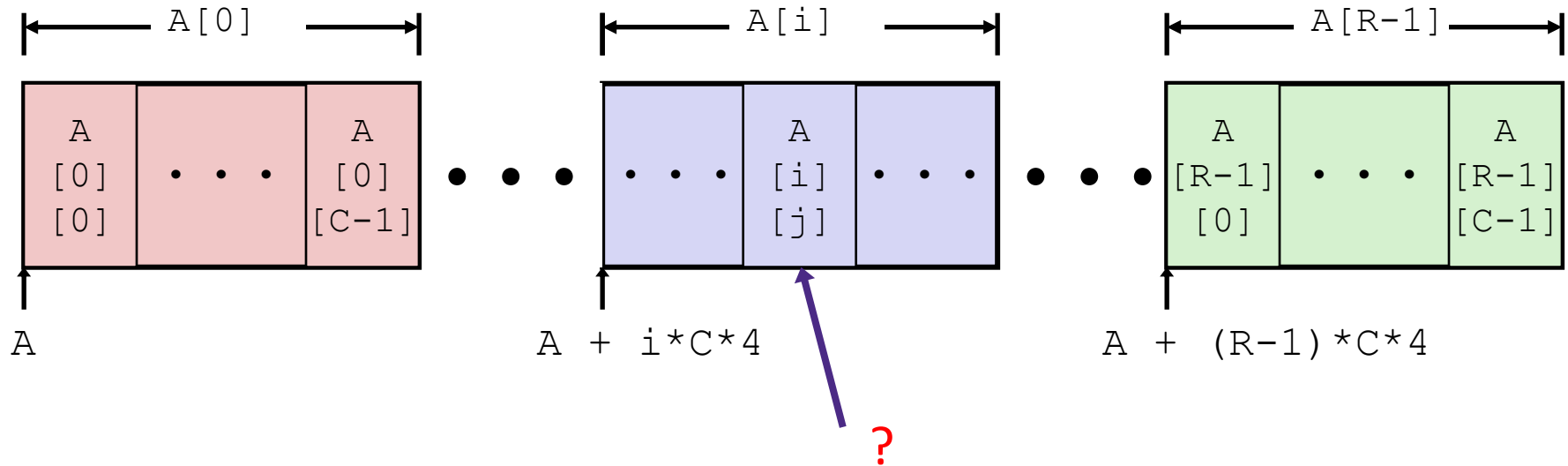  ▪ Starting address = `sea+20*index`

❖ Assembly Code
  ▪ Computes and returns address
  ▪ Compute as: `sea+4*(index+4*index)= sea+20*index`

34

# Nested Array <u>Element Access</u>

❖ Array Elements
  ▪ `A[i][j]` is element of type **T**, which requires $K$ bytes
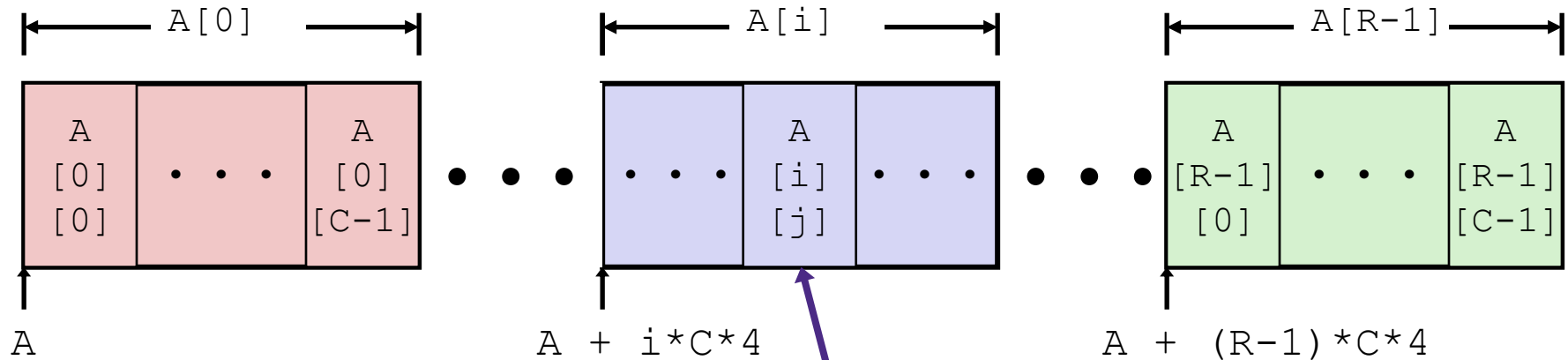  ▪ Address of `A[i][j]` is

**int** `A[R][C];`

# Nested Array <u>Element Access</u>

❖ Array Elements
  ▪ `A[i][j]` is element of type `T`, which requires $K$ bytes
  ▪ Address of `A[i][j]` is
$$A + i*(C*K) + j*K == A + (i*C + j)*K$$

`int A[R][C];`



$$A + i*C*4 + j*4$$

# Nested Array <u>Element Access</u> Code

```
int get_sea_digit
  (int index, int digit)
{
  return sea[index][digit];
}
```

```
int sea[4][5] =
  {{ 9, 8, 1, 9, 5 },
   { 9, 8, 1, 0, 5 },
   { 9, 8, 1, 0, 3 },
   { 9, 8, 1, 1, 5 }};
```

```
leaq    (%rdi,%rdi,4), %rax    # 5*index
addl    %rax, %rsi             # 5*index+digit
movl    sea(,%rsi,4),  %eax    # *(sea + 4*(5*index+digit))
```
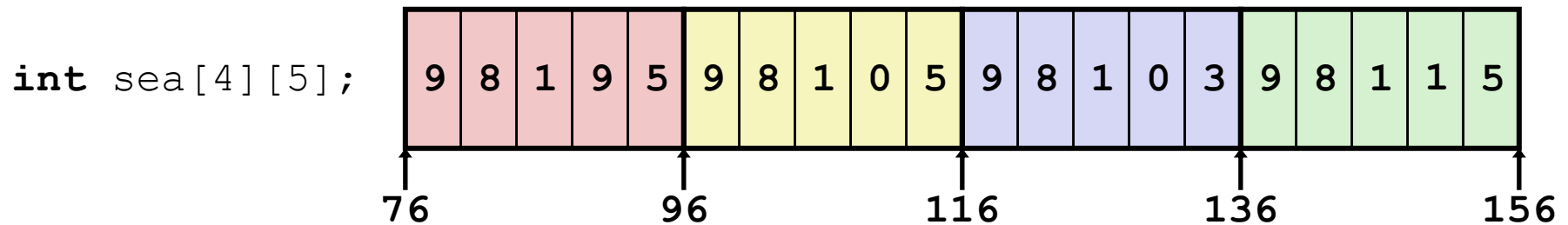
❖ Array Elements
- `sea[index][digit]` is an **int** (**sizeof**(**int**)=4)
- Address = `sea + 5*4*index + 4*digit`

❖ Assembly Code
- Computes address as:  `sea + ((index+4*index) + digit)*4`
- `movl` performs memory reference

# Multidimensional Referencing Examples

`int sea[4][5];`

| 9 | 8 | 1 | 9 | 5 | 9 | 8 | 1 | 0 | 5 | 9 | 8 | 1 | 0 | 3 | 9 | 8 | 1 | 1 | 5 |

76          96          116          136          156

## Reference  Address                          Value  Guaranteed?
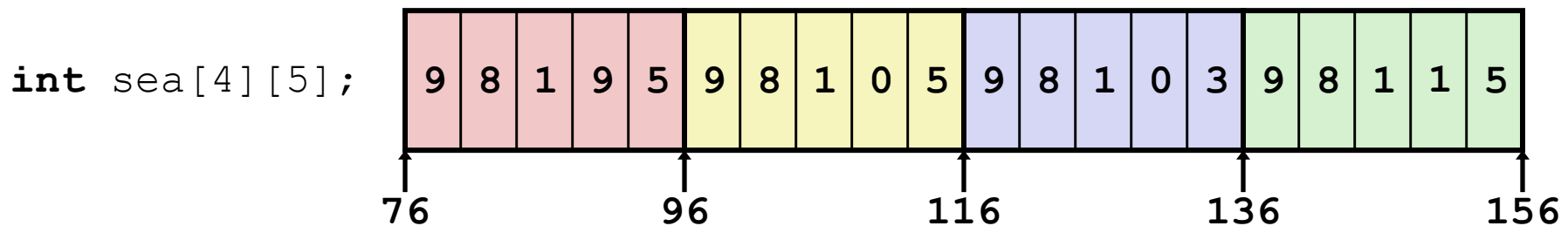
`sea[3][3]`

`sea[2][5]`

`sea[2][-1]`

`sea[4][-1]`

`sea[0][19]`

`sea[0][-1]`

- Code does not do any bounds checking
- Ordering of elements within array guaranteed

# Polling Question [Arrays - a]

❖ Which of the following statements is <u>FALSE</u>?

  ▪ Vote at http://pollev.com/pbjones

`int sea[4][5];`

| | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 9 | 8 | 1 | 9 | 5 | 9 | 8 | 1 | 0 | 5 | 9 | 8 | 1 | 0 | 3 | 9 | 8 | 1 | 1 | 5 |

76        96        116        136        156

A. `sea[4][-2]` is a *valid* array reference

B. `sea[1][1]` makes *two* memory accesses

C. `sea[2][1]` will *always* be a higher address than `sea[1][2]`

D. `sea[2]` is calculated using *only* `lea`

E. We're lost…

39

# Data Structures in Assembly

❖ **Arrays**
- One-dimensional
- Multidimensional (nested)
- **Multilevel**

❖ Structs
- Alignment

❖ ~~Unions~~

# **Multilevel** Array Example

**Multilevel Array Declaration(s):**

```
int cmu[5] = { 1, 5, 2, 1, 3 };
int  uw[5] = { 9, 8, 1, 9, 5 };
int ucb[5] = { 9, 4, 7, 2, 0 };
```

```
int* univ[3] = {uw, cmu, ucb};
```

Is a multilevel array the same thing as a 2D array?   **NO**

**2D Array Declaration:**

```
int univ2D[3][5] = {
  { 9, 8, 1, 9, 5 },
  { 1, 5, 2, 1, 3 },
  { 9, 4, 7, 2, 0 }
};
```
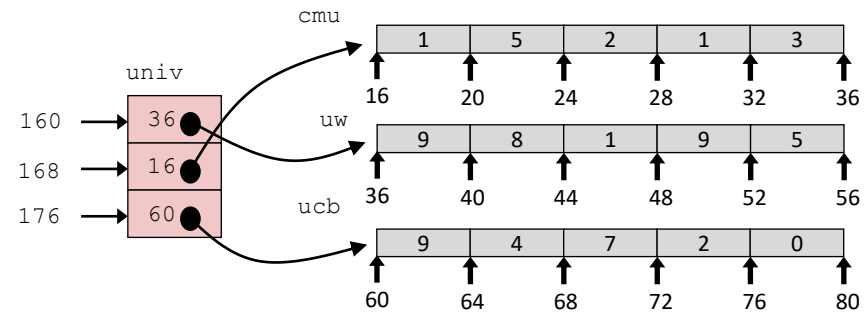
One array declaration = one contiguous block of memory

# **Multilevel Array Example**

```
int cmu[5] = { 1, 5, 2, 1, 3 };
int  uw[5] = { 9, 8, 1, 9, 5 };
int ucb[5] = { 9, 4, 7, 2, 0 };
```

```
int* univ[3] = {uw, cmu, ucb};
```

❖ Variable `univ` denotes array of 3 elements
❖ Each element is a pointer
  ▪ 8 bytes each
❖ Each pointer points to array of `int`s



Note: this is how Java represents multidimensional arrays

# Element Access in <u>Multilevel</u> Array

```
int get_univ_digit
   (int index, int digit)
{

   return univ[index][digit];

}
```



```
 salq     $2, %rsi             # rsi = 4*digit
 addq     univ(,%rdi,8), %rsi  # p = univ[index] + 4*digit
 movl     (%rsi), %eax         # return *p
 ret
```
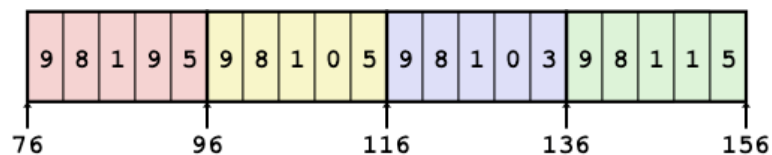
❖ Computation
  ▪ Element access  Mem[Mem[univ+8*index]+4*digit]
  ▪ Must do **two memory reads**
    • First get pointer to row array
    • Then access element within array
  ▪ But allows inner arrays to be different lengths (not in this example)

43

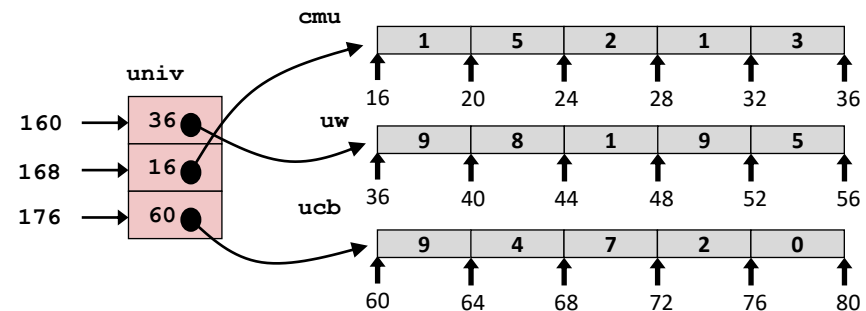# Array Element Accesses

## Multidimensional array

```
int get_sea_digit
  (int index, int digit)
{
  return sea[index][digit];
}
```

## Multilevel array

```
int get_univ_digit
  (int index, int digit)
{
  return univ[index][digit];
}
```
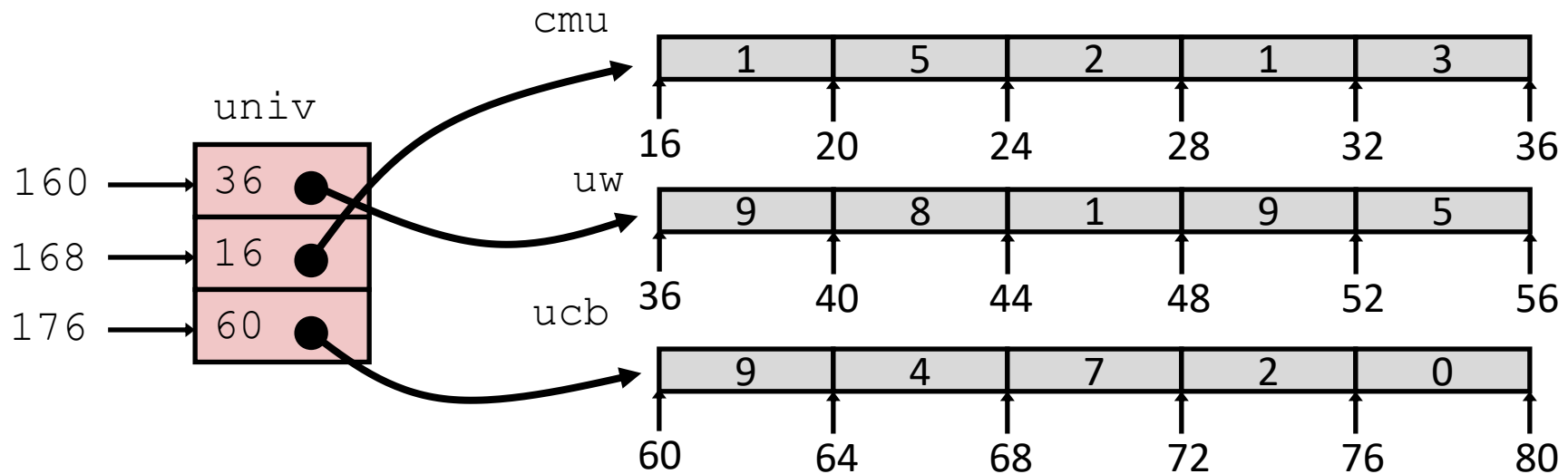


Access *looks* the same, but it isn't:

Mem[sea+20*index+4*digit]      Mem[Mem[univ+8*index]+4*digit]

# Multilevel Referencing Examples

cmu

| 1 | 5 | 2 | 1 | 3 |

16        20        24        28        32        36

univ

uw

| 9 | 8 | 1 | 9 | 5 |

160 → 36 ●

168 → 16 ●                               36        40        44        48        52        56

176 → 60 ●                     ucb

| 9 | 4 | 7 | 2 | 0 |

60        64        68        72        76        80

| Reference | Address | Value | Guaranteed? |
|-----------|---------|-------|-------------|
| univ[2][3] | | | |
| univ[1][5] | | | |
| univ[2][-2] | | | |
| univ[3][-1] | | | |
| univ[1][12] | | | |

- C code does not do any bounds checking
- Location of each lower-level array in memory is *not* guaranteed

45

# Summary

❖ Contiguous allocations of memory

❖ No bounds checking (and no default initialization)

❖ Can usually be treated like a pointer to first element

❖ `int a[4][5];` → array of arrays

- all levels in one contiguous block of memory

❖ `int* b[4];` → array of pointers to arrays

- First level in one contiguous block of memory

- Each element in the first level points to another "sub" array

- Parts anywhere in memory