

# Programming Design

## Pointers

Ling-Chieh Kung

Department of Information Management  
National Taiwan University

# Outline

- **Basics of pointers**
- Call by reference/pointer
- Arrays and pointer arithmetic
- Dynamic memory allocation (DMA)

# Pointers

- A **pointer** is a variable which stores a **memory address**.
  - An **array** variable is a pointer.
- To declare a pointer, use **\***.

```
type pointed* pointer name;
```

```
type pointed *pointer name;
```

- Examples:

```
int *ptrInt;
```

```
double* ptrDou;
```

- These pointers will store addresses.
  - These pointers will store addresses of **int/double** variables.
- We may point to **any** type.
- To point to different types, use different types of pointers.

# Sizes of pointers

- All pointers have the same size.
  - In a 32-bit computer, a pointer is allocated 4 bytes.
  - In a 64-bit computer, a pointer is allocated 8 bytes.

```
int* p1 = 0;  
cout << sizeof(p1) << endl; // 8  
double* p2 = 0;  
cout << sizeof(p2) << endl; // 8
```

- The length of pointers decides the maximum size of the memory space.
  - 32 bits:  $2^{32}$  bytes = 4GB.
  - 64 bits:  $2^{64}$  bytes = ?

# Pointer assignment

- We use the **address-of operator &** to obtain a variable's address:

```
pointer name = &variable name
```

- The address-of operator **&** returns the (beginning) **address** of a variable.
- Example:

- **ptr** points to **a**, i.e., **ptr** stores **the address of a**.

```
int a = 5;  
int* ptr = &a;
```

- When assigning an address, the two types must **match**.

```
int a = 5;  
double* ptr = &a; // error!
```

# Variables in memory

- `int a = 5;`
- `double b = 10.5;`
- `int* aPtr = &a;`
- `double* bPtr = &b;`
- `cout << &a; // 0x20c644`
- `cout << &b; // 0x20c660`
- `cout << &aPtr; // 0x20c658`
- `cout << &bPtr; // 0x20c64c`

Address	Identifier	Value
0x20c644	a	5
0x20c64c	bPtr	0x20c660
0x20c650		
0x20c658	aPtr	0x20c644
0x20c65c		
0x20c660	b	10.5
0x20c664		

Memory

# Address operators

- There are two address operators.
  - **&**: The **address-of operator**. It returns a variable's address.
  - **\***: The **dereference operator**. It returns the pointed variable (not the value!).
- For `int a = 5`:
  - `a` equals 5.
  - `&a` returns an address (e.g., 0x22ff78).
- For `int* ptrA = &a`:
  - `ptrA` stores an address (e.g., 0x22ff78).
  - `&ptrA` returns the pointer's address (e.g., 0x21aa74). This has nothing to do with the pointed variable `a`.
  - `*ptrA` returns `a`, **the variable** pointed by the pointer.

# Address operators

- Example:

```
int a = 10;
int* p1 = &a;
cout << "value of a = " << a << endl;
cout << "value of p1 = " << p1 << endl;
cout << "address of a = " << &a << endl;
cout << "address of p1 = " << &p1 << endl;
cout << "value of the variable pointed by p1 = " << *p1 << endl;
```

# Address operators and NULL

- **&**: returns a variable's address.
  - We cannot use **&100**, **&(a++)** (because **a++** returns the value of **a**).
  - We can only perform **&** on a **variable**.
  - We cannot assign value to **&x** (**&x** is a value!).
  - We can get a usual variable's or a pointer variable's address.
- **\***: returns the pointed variable, **not** its value.
  - We can perform **\*** on a pointer variable.
  - We cannot perform **\*** on a usual variable.
  - We cannot change a variable's address. No operation can do this.
- A pointer pointing to nothing should be assigned **NULL** or **0**.

# Address operators and NULL

- Examples:

```
int a = 10;
int* ptr = NULL;
ptr = &a;
cout << *ptr; // ?
*ptr = 5;
cout << a;    // ?
a = 18;
cout << *ptr; // ?
```

```
int a = 10;
int* ptr1 = NULL;
int* ptr2 = NULL;
ptr1 = ptr2 = &a;
cout << *ptr1; // ?
*ptr2 = 5;
cout << *ptr1; // ?
(*ptr1)++;
cout << a;    // ?
```

# Address operators and NULL

- Dereferencing a null pointer shutdowns the program (a run-time error).

```
int* p2 = NULL;
cout << "value of p2 = " << p2 << endl;
cout << "address of p2 = " << &p2 << endl;
cout << "the variable pointed by p2 = " << *p2 << endl;
```

# & and \* cancel each other

- What is  $*\&x$  if  $x$  is a variable?
  - $\&x$  is the address of  $x$ .
  - $*(\&x)$  is the variable stored in that address.
  - So  $*(\&x)$  is  $x$ .
- What is  $\&*x$  if  $x$  is a pointer?
  - If  $x$  is a pointer,  $*x$  is the variable stored at  $x$  ( $x$  stores an address!).
  - $\&*x$  is the address of  $*x$ , which is exactly  $x$ .
- What is  $\&*x$  if  $x$  is not a pointer?

# Good programming style

- Initialize a pointer variable as **0** or **NULL** if no initial value is available.
  - **0** is the standard in C++, while **NULL** is the standard in C. But they are the same for representing “null pointer”.
  - By using **NULL**, everyone knows the variable must be a pointer, and you are not talking about a number or character.
- Without an initialization, a pointer points to **somewhere**... And we do not know where it is!
  - Accessing an unknown address results in unpredictable results.
- In general, when you get **a run time error** or different outcomes for multiple executions, check your arrays and pointers.

# Good programming style

- As a bad example:

```
#include <iostream>
using namespace std;

int main()
{
    int* ptrArray[10000];
    for(int i = 0; i < 10000; i++)
        cout << i << " " << *ptrArray[i] << "\n";
    return 0;
}
```

# Good programming style

- When we use `*` in **declaring** a pointer, that `*` is not a dereference operator.
  - It is just a special syntax for declaring a pointer variable.
- I prefer to view `int*` as a type, which represents an “integer pointer”.
- Therefore, I prefer “`int* p`” to “`int *p`”.
- Be careful:

```
int* p, q;    // p is int*, q is int
int *p, *q;   // two pointers
int* p, *q;   // two pointers
int* p, * q;  // two pointers
```

- I use multiple statements to declare multiple pointers.

# Outline

- The basics of pointers
- **Call by reference/pointer**
- Arrays and pointer arithmetic
- Dynamic memory allocation (DMA)

# References and pointers

- Recall this example:

```
void swap (int x, int y);
int main()
{
    int a = 10, b = 20;
    cout << a << " " << b << endl;
    swap(a, b);
    cout << a << " " << b << endl;
}
void swap (int x, int y)
{
    int temp = x;
    x = y;
    y = temp;
}
```

# References and pointers

- When invoking a function and passing parameters, the default scheme is to “**call by value**” (or “pass by value”).
  - The function declares its own local variables, using a copy of the arguments’ values as initial values.
  - Thus we swapped the two local variables declared in the function, not the original two we want to swap.
- To solve this, we can use “**call by reference**” or “call by pointer.”
  - They are somewhat different, but the principle is the same.
  - It is enough to know and use only one of them.

# Call by reference

- A **reference** is a variable's alias.
- The reference is another variable that refers to the variable.
- Thus, using the reference is the same as using the variable.

```
int c = 10;  
int& d = c; // declare d as c's reference  
d = 20;  
cout << c << endl; // 20
```

- **int& d = c** is to declare **d** as **c**'s reference.
  - This **&** is different from the **&** operator which returns a variable's address.
- **int& d = 10** is an error.
  - A literal cannot have an alias!

# Call by reference

- Now we know how to change a parameter's value:
  - Instead of declaring a usual local variable as a parameter, declare a **reference** variable.
- This is to “call by reference”.

```
void swap (int& x, int& y);  
int main()  
{  
    int a = 10, b = 20;  
    cout << a << " " << b << endl;  
    swap(a, b);  
    cout << a << " " << b << endl;  
}  
void swap (int& x, int& y)  
{  
    int temp = x;  
    x = y;  
    y = temp;  
}
```

# Call by reference

- Thus we can call by reference and modify our parameters' value.
- When calling by reference, the only thing you have to do is to add an **&** in the parameter declaration in the function header.
- Mostly people use references only to call by reference.
- View the **&** in declaration as a part of type.
  - I use **int& a = b;** instead of **int &a = b;**.

```
void swap (int& x, int& y);  
int main()  
{  
    int a = 10, b = 20;  
    swap(a, b);  
}
```

# Call by pointers

- To call by pointers:
  - Declare a **pointer** variable as a parameter.
  - Pass a pointer variable or an address (returned by **&**) at invocation.
- For the **swap()** example:

```
void swap(int* ptrA, int* ptrB)
{
    int temp = *ptrA;
    *ptrA = *ptrB;
    *ptrB = temp;
}
```

- Invocation becomes **swap(&a, &b);**

# Call by pointers

- How about the following implementation?

```
void swap(int* ptrA, int* ptrB)
{
    int* temp = ptrA;
    ptrA = ptrB;
    ptrB = temp;
}
```

- Invocation: `swap(&a, &b);`
- Will the two arguments be swapped? What really happens?

# Call by pointers

- The principle behind calling by reference and calling by pointer is the same.
- You can view calling by reference as a special tool made by using pointers.
- Do not mix references and pointers!
  - E.g., we cannot pass a pointer variable or an address to a reference!
- You can use calling by reference in most situations, and it is clearer and more convenient than calling by pointer.
  - When you just want to modify arguments or return several values, call by reference.
  - When you really have to do something by pointers, call by pointer.

# Outline

- The basics of pointers
- Call by reference/pointer
- **Arrays and pointer arithmetic**
- Dynamic memory allocation (DMA)

# Pointers and arrays

- An array variable **is** a pointer!
  - It records the address of the **first** element of the array.
  - When passing an array, we pass a pointer.
  - The array indexing operator **[ ]** indicates **offsetting**.
- To further understand this issue, let's study **pointer arithmetic**.
  - Using **+**, **-**, **++**, and **--** on pointers.

# Pointer arithmetic

- Usually, one arbitrary address returned by performing arithmetic on a pointer variable is useless.
- The arithmetic is useful (and should be used) only when you can predict a variable's address.
  - In particular, when variables are stored **consecutively**.

```
int a = 10;
int* ptr = &a;
cout << ptr++;
    // just an address
    // we don't know what's here
cout << *ptr;
    // dangerous!
```

# Pointer Arithmetic: ++ and --

- **++**: Increment the pointer variable's value by the number of bytes occupied by a variable in this type (i.e., point to the **next** variable).
  - E.g., for integer pointers, the value (an address) increases by 4 (bytes).
- **--**: Decrement the pointer variable's value by the number of bytes a variable in this type occupies (i.e., point to the **previous** variable).

```
double a[3] = {10.5, 11.5, 12.5};
double* b = &a[0];
cout << *b << " " << b << endl; // 10.5
b = b + 2;
cout << *b << " " << b << endl; // 12.5
b--;
cout << *b << " " << b << endl; // 11.5
```

# Pointer Arithmetic: –

- We cannot add two address.
- However, we can find the difference of two addresses.

```
double a[3] = {10.5, 11.5, 12.5};  
double* b = &a[0];  
double* c = &a[2];  
cout << c - b << endl; // 2, not 16!
```

# Pointers and arrays

- Changing the value stored in a pointer is dangerous:

```
int y[3] = {1, 2, 3};
int* x = y;
for(int i = 0; i < 3; i++)
    cout << *(x + i) << " "; // 1 2 3
for(int i = 0; i < 3; i++)
    cout << *(x++) << " "; // 1 2 3
for(int i = 0; i < 3; i++)
    cout << *(x + i) << " "; // unpredictable
```

# Indexing and pointer arithmetic

- The array indexing operator `[]` is just an **interface** for doing pointer arithmetic.

```
int x[3] = {1, 2, 3};  
for(int i = 0; i < 3; i++)  
    cout << x[i] << " "; // x[i] = *(x + i)  
for(int i = 0; i < 3; i++)  
    cout << *(x++) << " "; // error!
```

- An array variable (e.g., `x`) stores an address, but `++` and `--` work only on pointer variables.
- Interface: a (typically safer and easier) way of completing a task.
  - `x[i]` and `*(x + i)` are identical.
  - But using the former is safer and easier.

# Example: insertion sort

- Consider the **insertion sort** taught last time.
  - Given a unsorted array  $A$  of length  $n$ , we first sort  $A[0:(n-2)]$ , and then insert  $A[n-1]$  to the sorted part.
  - To complete this task, we do this **recursively**.
- What if we want to **first sort  $A[1:(n-1)]$** , and then insert  $A[0]$ ?
- We will need to implement a function:
  - **`void insertionSort(int array[], const int n);`**
  - Given **`array`**, each time when we (recursively) invoke it, we pass a shorter array formed by elements from **`array[1]`** to **`array[n - 1]`**.
  - How?

# Example: insertion sort

```
void insertionSort(int array[], const int n) {
    if(n > 1) {
        insertionSort(array + 1, n - 1);
        int num1 = array[0];
        int i = 1;
        for(; i < n; i++) {
            if(array[i] < num1)
                array[i - 1] = array[i];
            else
                break;
        }
        array[i - 1] = num1;
    }
}
```

# Outline

- The basics of pointers
- Call by reference/pointer
- Arrays and pointer arithmetic
- **Dynamic memory allocation (DMA)**

# Static memory allocation

- In C/C++, we declare an array by specifying its length as a constant variable or a literal.
  - `int a[100];`
- A memory space will be allocated to an array during the compilation time.
  - 400 bytes will be allocated for the above statement.
- This is called “**static memory allocation**”.
- We may decide the length of an array “**dynamically**”.
  - That is, during the **run** time.
- To do so, we must use a different syntax.
  - All types of variables may also be declared in this way.

# Dynamic memory allocation

- The operator **new** allocates a memory space **and** returns the address.
  - In C, we use a different keyword **malloc**.
- **new int**; allocates 4 bytes without recording the address.
- **int\* a = new int**; makes **a** store the address of the space.
- **int\* a = new int(5)** ; makes the space contains 5 as the value.
- **int\* a = new int[5]** ; allocates 20 bytes (for 5 integers).
  - **a** points to the first integer.
- Dynamically allocated arrays **cannot be initialized** with a single statement.
  - A loop, for example, is needed.

# Dynamic memory allocation

- All of these spaces are allocated during the **run time**.
- So we may write

```
int len = 0;  
cin >> len;  
int* a = new int[len];
```

- This allocates a space according to the input from users.

# Dynamic memory allocation

- A space allocated during the run time has **no name!**
  - On the other hand, every space allocated during compilation time has a name.
- To access a dynamically-allocated space, we use a **pointer** to store its address.

```
int len = 0;
cin >> len; // 3
int* a = new int[len];
for (int i = 0; i < len; i++)
    a[i] = i + 1;
```

Address	Identifier	Value
0x20c644	N/A	1
0x20c648		2
0x20c64c		3
0x20c650		
0x20c654		
0x20c658	len	3
0x20c65c		
0x20c660	a	0x20c644
0x20c664		

Memory

# Example: Fibonacci sequence

- Recall the repetitive implementation of generating the Fibonacci sequence.
- After we get the value of sequence length  $n$ , we dynamically declare an array of length  $n$ .
- Then just use that array!

```
double fibRepetitive (int n)
{
    if (n == 1)
        return 1;
    else if (n == 2)
        return 1;
    double* fib = new double[n];
    fib[0] = 1;
    fib[1] = 1;
    for (int i = 2; i < n; i++)
        fib[i] = fib[i - 1] + fib[i - 2];
    double result = fib[n - 1];
    delete[] fib; // to be explained
    return result;
}
```

# Memory leak

- For spaces allocated during the **compilation** time, the system will **release these spaces** automatically when the corresponding variables no longer exist.
- For spaces allocated during the **run** time, the system will **NOT** release these spaces unless it is asked to do so.
  - Because the space has no name!

```
void func (int a)
{
    double b;
} // 4 + 8 bytes are released
```

```
void func()
{
    int* bPtr = new int[10];
}
// 8 bytes for bPtr are released
// 40 bytes for integers are not
```

# Memory leak

- Programmers must keep a record for all spaced allocated dynamically.

```
double* b = new double;  
*b = 5.2;  
double c = 10.6;  
b = &c; // now no one can access  
        // the space containing 5.2
```

- This problem is called **memory leak**.
  - We lose the control of allocated spaces.
  - These spaces are **wasted**.
  - They will not be released until the program ends.

# Memory leak

- Try this carefully!
  - The outcome may be different on your computer.

```
#include <iostream>
using namespace std;

int main()
{
    for(int i = 0; ; i++)
    {
        int* ptr = new int[10000];
        cout << i << "\n";
        // delete [] ptr;
    }
    return 0;
}
```

# Releasing spaces manually

- The **delete** operator will release a dynamically-allocated space.

```
int* a = new int;  
delete a; // release 4 bytes  
int* b = new int[5];  
delete b; // release only 4 bytes!  
           // Unpredictable results may happen  
delete [] b; // release all 20 bytes
```

- The **delete** operator will do nothing to the pointer. To avoid reuse the released space, set the pointer to **NULL**.

```
int* a = new int;  
delete a; // a is still pointing to the address  
a = NULL; // now a points to nothing  
int* b = new int[5];  
delete [] b; // b is still pointing to the address  
b = NULL; // now b points to nothing
```

# Good programming style

- Use DMA for arrays with **no predetermined** length.
  - Even though Dev-C++ (and many other compilers) converts

```
int a = 10;  
int b[a];
```

to

```
int a = 10;  
int* b = new int[a];
```

- To avoid memory leak:
  - Whenever you write a **new** statement, add a **delete** statement below immediately (unless you know you really do not need it).
  - Whenever you want to change the value of a pointer, check whether memory leak occurs.
  - Whenever you write a **delete** statement, set the pointer to **NULL**.

# Two-dimensional dynamic arrays

- With static arrays, we may create matrices as two-dimensional arrays.
- An  $m$  by  $n$  two-dimensional array has:
  - $m$  rows (single-dimensional arrays).
  - Each row has  $n$  elements.
- With dynamic arrays, we now may create matrices **with different row lengths**.
  - We may still have  $m$  rows.
  - Now each row may have different number of elements.
  - E.g., a **lower triangular matrix**.

# Example: lower triangular arrays

```
#include <iostream>
using namespace std;

int print(int** arr, int r)
{
    for(int i = 0; i < r; i++)
    {
        for(int j = 0; j < i; j++)
            cout << arr[i][j] << " ";
        cout << "\n";
    }
}
```

```
int main()
{
    int r = 10;
    int** array = new int*[r];
    for(int i = 0; i < r; i++)
    {
        array[i] = new int[i + 1];
        for(int j = 0; j <= i; j++)
            array[i][j] = j + 1;
    }
    print(array, r);
    return 0;
}
```

# A pointer for pointers

- `int* array = new int[10];` declares **a pointer for integers**.
- `int** array = new int*[10];` declares **a pointer for integer pointers!**
  - The type of `array[0]` is `int*`.
  - The type of `array[1]` is `int*`.
- Then each of these `int*` may point to one or multiple integers.
  - And their lengths can be different.

```
int r = 10;
int** array = new int*[r];
for(int i = 0; i < r; i++)
    array[i] = new int[i + 1];
```