

M.1 — Intro to smart pointers and move semantics

ALEX JANUARY 25, 2022

Consider a function in which we dynamically allocate a value:

```
1 void someFunction()
2 {
3     Resource *ptr = new Resource(); // Resource is a struct or class
4
5     // do stuff with ptr here
6
7     delete ptr;
8 }
```

Although the above code seems fairly straightforward, it's fairly easy to forget to deallocate ptr. Even if you do remember to delete ptr at the end of the function, there are a myriad of ways that ptr may not be deleted if the function exits early. This can happen via an early return:

```
1 #include <iostream>
2
3 void someFunction()
4 {
5     Resource *ptr = new Resource();
6
7     int x;
8     std::cout << "Enter an integer: ";
9     std::cin >> x;
10
11     if (x == 0)
12         return; // the function returns early, and ptr won't be deleted!
13
14     // do stuff with ptr here
15
16     delete ptr;
17 }
```

or via a thrown exception:

```
1 #include <iostream>
2
3 void someFunction()
4 {
5     Resource *ptr = new Resource();
6
7     int x;
8     std::cout << "Enter an integer: ";
9     std::cin >> x;
10
11     if (x == 0)
12         throw 0; // the function returns early, and ptr won't be deleted!
13
14     // do stuff with ptr here
15
16     delete ptr;
17 }
```

In the above two programs, the early return or throw statement execute, causing the function to terminate without variable ptr being deleted. Consequently, the memory allocated for variable ptr is now leaked (and will be leaked again every time this function is called and returns early).

At heart, these kinds of issues occur because pointer variables have no inherent mechanism to clean up after themselves.

Smart pointer classes to the rescue?

One of the best things about classes is that they contain destructors that automatically get executed when an object of the class goes out of scope. So if you allocate (or acquire) memory in your constructor, you can deallocate it in your destructor, and be guaranteed that the memory will be deallocated when the class object is destroyed (regardless of whether it goes out of scope, gets explicitly deleted, etc...). This is at the heart of the RAII programming paradigm that we talked about in lesson 12.9 -- Destructors.

So can we use a class to help us manage and clean up our pointers? We can!

Consider a class whose sole job was to hold and "own" a pointer passed to it, and then deallocate that pointer when the class object went out of scope. As long as objects of that class were only created as local variables, we could guarantee that the class would properly go out of scope (regardless of when or how our functions terminate) and the owned pointer would get destroyed.

Here's a first draft of the idea:

```
1 #include <iostream>
2
3 template <typename T>
4 class Auto_ptr1
5 {
6     T* m_ptr;
7 public:
8     // Pass in a pointer to "own" via the constructor
9     Auto_ptr1(T* ptr=nullptr)
10     {
11         :m_ptr(ptr)
12     }
13
14     // The destructor will make sure it gets deallocated
15     ~Auto_ptr1()
16     {
17         delete m_ptr;
18     }
19
20     // Overload dereference and operator-> so we can use Auto_ptr1 like m_ptr.
21     T& operator*() const { return *m_ptr; }
22     T* operator->() const { return m_ptr; }
23 };
24
25 // A sample class to prove the above works
26 class Resource
27 {
28 public:
29     Resource() { std::cout << "Resource acquired\n"; }
30     ~Resource() { std::cout << "Resource destroyed\n"; }
31 };
32
33 int main()
34 {
35     Auto_ptr1<Resource> res(new Resource()); // Note the allocation of memory here
36
37     // ... but no explicit delete needed
38
39     // Also note that the Resource in angled braces doesn't need a * symbol, since that's supplied by the template
40
41     return 0;
42 } // res goes out of scope here, and destroys the allocated Resource for us
```

This program prints:

```
Resource acquired
Resource destroyed
```

Consider how this program and class work. First, we dynamically create a Resource, and pass it as a parameter to our templated Auto_ptr1 class. From that point forward, our Auto_ptr1 variable res owns that Resource object (Auto_ptr1 has a composition relationship with m_ptr). Because res is declared as a local variable and has block scope, it will go out of scope when the block ends, and be destroyed (no worries about forgetting to deallocate it). And because it is a class, when it is destroyed, the Auto_ptr1 destructor will be called. That destructor will ensure that the Resource pointer it is holding gets deallocated!

As long as Auto_ptr1 is defined as a local variable (with automatic duration, hence the "Auto" part of the class name), the Resource will be guaranteed to be destroyed at the end of the block it is declared in, regardless of how the function terminates (even if it terminates early).

Such a class is called a smart pointer. A **Smart pointer** is a composition class that is designed to manage dynamically allocated memory and ensure that memory gets deleted when the smart pointer object goes out of scope. (Relatedly, built-in pointers are sometimes called "dumb pointers" because they can't clean up after themselves).

Now let's go back to our someFunction() example above, and show how a smart pointer class can solve our challenge:

```
1 #include <iostream>
2
3 template <typename T>
4 class Auto_ptr1
5 {
6     T* m_ptr;
7 public:
8     // Pass in a pointer to "own" via the constructor
9     Auto_ptr1(T* ptr=nullptr)
10     {
11         :m_ptr(ptr)
12     }
13
14     // The destructor will make sure it gets deallocated
15     ~Auto_ptr1()
16     {
17         delete m_ptr;
18     }
19
20     // Overload dereference and operator-> so we can use Auto_ptr1 like m_ptr.
21     T& operator*() const { return *m_ptr; }
22     T* operator->() const { return m_ptr; }
23 };
24
25 // A sample class to prove the above works
26 class Resource
27 {
28 public:
29     Resource() { std::cout << "Resource acquired\n"; }
30     ~Resource() { std::cout << "Resource destroyed\n"; }
31     void sayHi() { std::cout << "Hi!\n"; }
32 };
33
34 void someFunction()
35 {
36     Auto_ptr1<Resource> ptr(new Resource()); // ptr now owns the Resource
37
38     int x;
39     std::cout << "Enter an integer: ";
40     std::cin >> x;
41
42     if (x == 0)
43         return; // the function returns early
44
45     // do stuff with ptr here
46     ptr->sayHi();
47 }
48
49 int main()
50 {
51     someFunction();
52
53     return 0;
54 }
```

If the user enters a non-zero integer, the above program will print:

```
Resource acquired
Hi!
Resource destroyed
```

If the user enters zero, the above program will terminate early, printing:

```
Resource acquired
Resource destroyed
```

Note that even in the case where the user enters zero and the function terminates early, the Resource is still properly deallocated.

Because the ptr variable is a local variable, ptr will be destroyed when the function terminates (regardless of how it terminates). And because the Auto_ptr1 destructor will clean up the Resource, we are assured that the Resource will be properly cleaned up.

A critical flaw

The Auto_ptr1 class has a critical flaw lurking behind some auto-generated code. Before reading further, see if you can identify what it is. We'll wait...

(Hint: consider what parts of a class get auto-generated if you don't supply them)

(Jeopardy music)

Okay, time's up.

Rather than tell you, we'll show you. Consider the following program:

```
1 #include <iostream>
2
3 // Same as above
4 template <typename T>
5 class Auto_ptr1
6 {
7     T* m_ptr;
8 public:
9     Auto_ptr1(T* ptr=nullptr)
10     {
11         :m_ptr(ptr)
12     }
13
14     ~Auto_ptr1()
15     {
16         delete m_ptr;
17     }
18
19     T& operator*() const { return *m_ptr; }
20     T* operator->() const { return m_ptr; }
21 };
22
23 class Resource
24 {
25 public:
26     Resource() { std::cout << "Resource acquired\n"; }
27     ~Resource() { std::cout << "Resource destroyed\n"; }
28 };
29
30 int main()
31 {
32     Auto_ptr1<Resource> res1(new Resource());
33     Auto_ptr1<Resource> res2(res1); // Alternatively, don't initialize res2 and then assign res2 = res1;
34     res1;
35
36     return 0;
37 }
```

This program prints:

```
Resource acquired
Resource destroyed
Resource destroyed
```

Very likely (but not necessarily) your program will crash at this point. See the problem now? Because we haven't supplied a copy constructor or an assignment operator, C++ provides one for us. And the functions it provides do shallow copies. So when we initialize res2 with res1, both Auto_ptr1 variables are pointed at the same Resource. When res2 goes out of the scope, it deletes the resource, leaving res1 with a dangling pointer. When res1 is deleted, both its (already deleted) Resource, crash!

You'd run into a similar problem with a function like this:

```
1 void passByValue(Auto_ptr1<Resource> res)
2 {
3 }
4
5 int main()
6 {
7     Auto_ptr1<Resource> res1(new
8     Resource());
9     passByValue(res1)
10
11     return 0;
12 }
```

In this program, res1 will be copied by value into passByValue's parameter res, leading to duplication of the Resource pointer. Crash!

So clearly this isn't good. How can we address this?

Well, one thing we could do would be to explicitly define and delete the copy constructor and assignment operator, thereby preventing any copies from being made in the first place. That would prevent the pass by value case (which is good, we probably shouldn't be passing these by value anyway).

But then how would we return an Auto_ptr1 from a function back to the caller?

```
1 ??? generateResource()
2 {
3     Resource* r{ new Resource() };
4     return Auto_ptr1(r);
5 }
```

We can't return our Auto_ptr1 by reference, because the local Auto_ptr1 will be destroyed at the end of the function, and the caller will be left with a dangling reference. We could return pointer r as Resource*, but then we might forget to delete r later, which is the whole point of using smart pointers in the first place. So that's out. Returning the Auto_ptr1 by value is the only option that makes sense -- but then we end up with shallow copies, duplicated pointers, and crashes.

Another option would be to override the copy constructor and assignment operator to make deep copies. In this way, we'd at least guarantee to avoid duplicate pointers to the same object. But copying can be expensive (and may not be desirable or even possible), and we don't want to make needless copies of objects just to return an Auto_ptr1 from a function. Plus assigning or initializing a dumb pointer doesn't copy the object being pointed to, so why would we expect smart pointers to behave differently?

What do we do?

Move semantics

What if, instead of having our copy constructor and assignment operator copy the pointer ("copy semantics"), we instead transfer/move ownership of the pointer from the source to the destination object? This is the core idea behind move semantics. **Move semantics** means the class will transfer ownership of the object rather than making a copy.

Let's update our Auto_ptr1 class to show how this can be done:

```
1 #include <iostream>
2
3 template <typename T>
4 class Auto_ptr2
5 {
6     T* m_ptr;
7 public:
8     Auto_ptr2(T* ptr=nullptr)
9     {
10         :m_ptr(ptr)
11     }
12
13     ~Auto_ptr2()
14     {
15         delete m_ptr;
16     }
17
18     // A copy constructor that implements move semantics
19     Auto_ptr2(Auto_ptr2& a) // note: not const
20     {
21         m_ptr = a.m_ptr; // transfer our dumb pointer from the source to our local object
22         a.m_ptr = nullptr; // make sure the source no longer owns the pointer
23     }
24
25     // An assignment operator that implements move semantics
26     Auto_ptr2& operator=(Auto_ptr2& a) // note: not const
27     {
28         if (&a == this)
29             return *this;
30
31         delete m_ptr; // make sure we deallocate any pointer the destination is already holding first
32         m_ptr = a.m_ptr; // then transfer our dumb pointer from the source to the local object
33         a.m_ptr = nullptr; // make sure the source no longer owns the pointer
34         return *this;
35     }
36
37     T& operator*() const { return *m_ptr; }
38     T* operator->() const { return m_ptr; }
39     bool isNull() const { return m_ptr == nullptr; }
40 };
41
42 class Resource
43 {
44 public:
45     Resource() { std::cout << "Resource acquired\n"; }
46     ~Resource() { std::cout << "Resource destroyed\n"; }
47 };
48
49 int main()
50 {
51     Auto_ptr2<Resource> res1(new Resource());
52     Auto_ptr2<Resource> res2; // Start as nullptr
53
54     std::cout << "res1 is " << (res1.isNull() ? "null\n" : "not null\n");
55     std::cout << "res2 is " << (res2.isNull() ? "null\n" : "not null\n");
56
57     res2 = res1; // res2 assumes ownership, res1 is set to null
58
59     std::cout << "Ownership transferred\n";
60
61     std::cout << "res1 is " << (res1.isNull() ? "null\n" : "not null\n");
62     std::cout << "res2 is " << (res2.isNull() ? "null\n" : "not null\n");
63
64     return 0;
65 }
```

This program prints:

```
Resource acquired
res1 is not null
res2 is null
Ownership transferred
res1 is null
res2 is not null
Resource destroyed
```

Note that our overloaded operator= gave ownership of m_ptr from res1 to res2! Consequently, we don't end up with duplicate copies of the pointer, and everything gets tidily cleaned up.

std:auto_ptr, and why it was a bad idea

Now would be an appropriate time to talk about std:auto_ptr. std:auto_ptr, introduced in C++98 and removed in C++17, was C++'s first attempt at a standardized smart pointer. std:auto_ptr opted to implement move semantics just like the Auto_ptr2 class does.

However, std:auto_ptr (and our Auto_ptr2 class) has a number of problems that makes using it dangerous.

First, because std:auto_ptr implements move semantics through the copy constructor and assignment operator, passing a std:auto_ptr by value to a function will cause your resource to get moved to the function parameter (and be destroyed at the end of the function when the function parameters go out of scope). Then when you go to access your auto_ptr argument from the caller (not realizing it was transferred and deleted), you're suddenly dereferencing a null pointer. Crash!

Second, std:auto_ptr always deletes its contents using non-array delete. This means auto_ptr won't work correctly with dynamically allocated arrays, because it uses the wrong kind of deallocation. Worse, it won't prevent you from passing it a dynamic array, which it will then mismanage, leading to memory leaks.

Finally, auto_ptr doesn't play nice with a lot of the other classes in the standard library, including most of the containers and algorithms. This occurs because those standard library classes assume that when they copy an item, it actually makes a copy, not a move.

Because of the above mentioned shortcomings, std:auto_ptr has been deprecated in C++11 and removed in C++17.

Moving forward

The core problem with the design of std:auto_ptr is that prior to C++11, the C++ language simply had no mechanism to differentiate "copy semantics" from "move semantics". Overriding the copy semantics to implement move semantics leads to weird edge cases and inadvertent bugs. For example, you can write `res1 = res2` and have no idea whether res2 will be changed or not!

Because of this, in C++11, the concept of "move" was formally defined, and "move semantics" were added to the language to properly differentiate copying from moving. Now that we've set the stage for why move semantics can be useful, we'll explore the topic of move semantics throughout the rest of this chapter. We'll also fix our Auto_ptr2 class using move semantics.

In C++11, std:auto_ptr has been replaced by a bunch of other types of "move-aware" smart pointers: std::unique_ptr, std::weak_ptr, and std::shared_ptr. We'll also explore the two most popular of these: unique_ptr (which is a direct replacement for auto_ptr) and shared_ptr.

Next lesson
M.2 R-value references

Back to table of contents

Previous lesson
20.x Chapter 20 comprehensive quiz