16.5 — Returning std::vector, and an introduction to move semantics

ALEX IANUARY 11, 2024

When we need to pass a std::vector to a function, we pass it by (const) reference so that we do not make an expensive copy of the array data.

Therefore, you will probably be surprised to find that it is okay to return a std::vector by value.

Copy semantics

Consider the following program:

```
#include <iostream>
#include <vector>

int main()
{
    std::vector arr1 { 1, 2, 3, 4, 5 }; // copies { 1, 2, 3, 4, 5 } into arr1
    std::vector arr2 { arr1 }; // copies arr1 into arr2

arr1[0] = 6; // We can continue to use arr1
    arr2[0] = 7; // and we can continue to use arr2

std::cout << arr1[0] << arr2[0] << '\n';
    return 0;
}</pre>
```

When arr2 is initialized with arr1, the copy constructor of std::vector is called, which copies arr1 into arr2.

Making a copy is the only reasonable thing to do in this case, as we need both arr1 and arr2 to live on independently. This example ends up making two copies, one for each initialization.

The term **copy semantics** refers to the rules that determine how copies of objects are made. When we say a type supports copy semantics, we mean that objects of that type are copyable, because the rules for making such copies have been defined. When we say copy semantics are being invoked, that means we've done something that will make a copy of an object.

. . .

0

For class types, copy semantics are typically implemented via the copy constructor (and copy assignment operator), which defines how objects of that type are copied. Typically this results in each data member of the class type being copied. In the prior example, the statement std::vector arr2 { arr1 }; invokes copy semantics, resulting in a call to the copy constructor of std::vector, which then makes a copy of each data member of arr1 into arr2. The end result is that arr1 is equivalent to (but independent from) arr2.

When copy semantics is not optimal

Now consider this related example:

```
#include <iostream>
#include <vector>

std::vector<int> generate() // return by value
{
    // We're intentionally using a named object here so mandatory copy elision doesn't apply
    std::vector arr1 { 1, 2, 3, 4, 5 }; // copies { 1, 2, 3, 4, 5 } into arr1
    return arr1;
}

int main()
{
    std::vector arr2 { generate() }; // the return value of generate() dies at the end of the expression
    // There is no way to use the return value of generate() here
    arr2[0] = 7; // we only have access to arr2

    std::cout << arr2[0] << '\n';
    return 0;
}</pre>
```

When arr2 is initialized this time, it is being initialized using a temporary object returned from function <code>generate()</code>. Unlike the prior case, where the initializer was an Ivalue that could be used in future statements, in this case, the temporary object is an rvalue will be destroyed at the end of the initialization expression. The temporary object can't be used beyond that point. Because the temporary (and its data) will be destroyed at the end of the expression, we need some way to get the data out of the temporary and into <code>arr2</code>.

The usual thing to do here is the same as in the previous example: use copy semantics and make a potentially expensive copy. That way arr2 gets its own copy of the data that can be used even after the temporary (and its data) is destroyed.

• • •

0

However, what makes this case different than the previous example is that the temporary is going to be destroyed anyway. After initialization is complete, the temporary doesn't need its data any more (which is why we can destroy it). We don't need two sets of data to exist simultaneously. In such cases, making a potentially expensive copy and then destroying the original data is suboptimal.

Introduction to move semantics

Instead, what if there was a way for arr2 to "steal" the temporary's data instead of copying it? arr2 would then be the new owner of the data, and no copy of the data would need to be made. When ownership of data is transferred from one object to another, we say that data has been **moved**. The cost of such a move is typically trivial (usually just two or three pointer assignments, which is way faster than copying an array of data!).

As an added benefit, when the temporary was then destroyed at the end of the expression, it would no longer have any data to destroy, so we wouldn't have to pay that cost either.

This is the essence of **move semantics**, which refers to the rules that determine how the data from one object is moved to another object. When move semantics is invoked, any data member that can be moved is moved, and any data member that can't be moved is copied. The ability to move data instead of copying it can make move semantics more efficient than copy semantics, especially when we can replace an expensive copy with an inexpensive move.

Key insight

Move semantics is an optimization that allows us, under certain circumstances, to inexpensively transfer ownership of some data members from one object to another object (rather than making a more expensive copy).

Data members that can't be moved are copied instead.

How move semantics is invoked

• •

0

Normally, when an object is being initialized with or assigned an object of the same type, copy semantics will be used (assuming the copy isn't elided).

However, when all of the following are true, move semantics will be invoked instead:

- The type of the object supports move semantics.
- The initializer or object being assigned from is an rvalue (temporary) object.
- The move isn't elided.

Here's the sad news: not that many types support move semantics. However, std::vector and std::string both do!

We'll dig into how move semantics works in more detail in <u>chapter 22</u> (https://www.learncpp.com#Chapter22). For now, it's enough to know what move semantics is, and which types are move-capable.

We can return move-capable types like std::vector by value

. . .

0

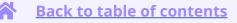
Because return by value returns an rvalue, if the returned type supports move semantics, then the returned value can be moved instead of copied into the destination object!

This makes return by value extremely inexpensive for these types!

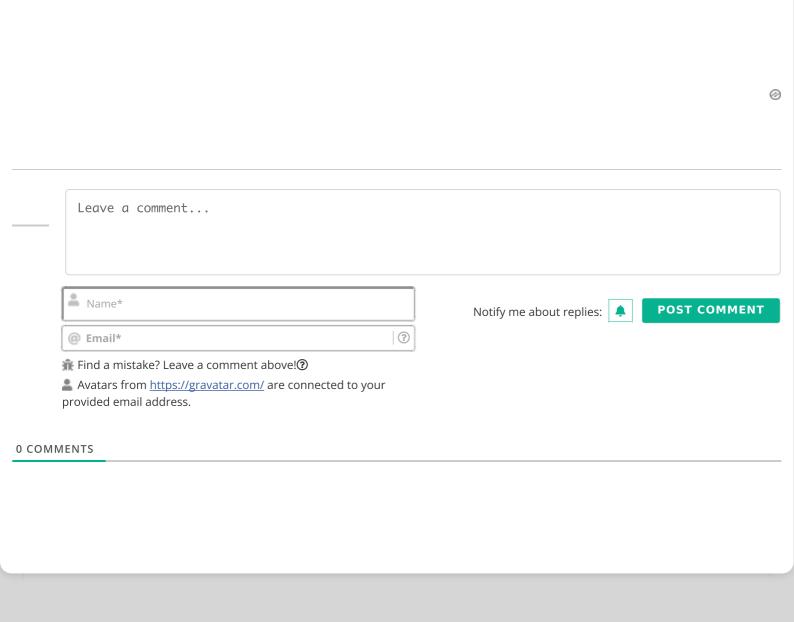
Key insight

We can return move-capable types (like std::string) by value. Such types will inexpensively move their values instead of making an expensive copy.









We and our partners share information on your use of this website to help improve your experience.

Do not sell my info:

OKAY

×