

# Move semantics and rvalue references in C++11



By Alex Allain

C++ has always produced fast programs. Unfortunately, until C++11, there has been an obstinate wart that slows down many C++ programs: the creation of temporary objects. Sometimes these temporary objects can be optimized away by the compiler (the return value optimization, for example). But this is not always the case, and it can result in expensive object copies. What do I mean?

Let's say that you have the following code:

```
1  #include <iostream>
2
3  using namespace std;
4
5  vector<int> doubleValues (const vector<
6  {
7      vector<int> new_values;
8      new_values.reserve(v.size());
9      for (auto itr = v.begin(), end_itr
10     {
11         new_values.push_back( 2 * *itr
12     }
13     return new_values;
14 }
15
16 int main()
17 {
18     vector<int> v;
19     for ( int i = 0; i < 100; i++ )
20     {
21         v.push_back( i );
22     }
23     v = doubleValues( v );
24 }
```

If you've done a lot of high performance work in C++, sorry about the pain that brought on. If you haven't--well, let's walk through why this code is terrible C++03 code. (The rest of this tutorial will be about why it's fine C++11 code.) The problem is with the copies. When `doubleValues` is called, it constructs a `vector`, `new_values`, and fills it up. This alone might not be ideal performance, but if we want to keep our original vector unsullied, we need a second copy. But what happens when we hit the return statement?

The entire contents of `new_values` must be copied! In principle, there could be up to two copies here: one into a temporary object to be returned, and a second when the vector assignment operator runs on the line `v = doubleValues( v );`. The first copy may be optimized away by the compiler automatically, but there is no avoiding that the assignment to `v` will have to copy all the values again, which requires a new memory allocation and another iteration over the entire vector.

This example might be a little bit contrived--and of course you can find ways to avoid this kind of problem--for example, by storing and returning the vector by pointer, or by passing in a vector to be filled up. The thing is, neither of these programming styles is particularly natural. Moreover, an approach that requires returning a pointer has introduced at least one more memory allocation, and one of the design goals of C++ is to avoid memory allocations.

The worst part of this whole story is that the object returned from `doubleValues` is a temporary value that's no longer needed. When you have the line `v = doubleValues( v );`, the result of `doubleValues( v )` is just going to get thrown away once it is copied! In theory, it should be possible to skip the whole copy and just pilfer the pointer inside the temporary vector and keep it in `v`. In effect, why can't we **move** the object? In C++03, the answer is that there was no way to tell if an object was a temporary or not, you had to run the same code in the assignment operator or copy constructor, no matter where the value came from, so no pilfering was possible. In C++11, the answer is--you can!

That's what rvalue references and move semantics are for! Move semantics allows you to avoid unnecessary copies when working with temporary objects that are about to evaporate, and whose resources can safely be taken from that temporary object and used by another.

Move semantics relies on a new feature of C++11, called rvalue references, which you'll want to understand to really appreciate what's going on. So first let's talk about what an rvalue is, and then what an rvalue reference is. Finally, we'll come back to move semantics and how it can be implemented with rvalue references.

## Rvalues and lvalues - bitter rivals, or best of friends?

In C++, there are rvalues and lvalues. An lvalue is an expression whose address can be taken, a locator value--essentially, an lvalue provides a (semi)permanent piece of memory. You can make assignments to lvalues. For example:

```
1 | int a;  
2 | a = 1; // here, a is an lvalue
```

You can also have lvalues that aren't variables:

```
1 | int x;  
2 | int& getRef ()  
3 | {  
4 |     return x;  
5 | }  
6 |  
7 | getRef() = 4;
```

Here, getRef returns a reference to a global variable, so it's returning a value that is stored in a permanent location. (You could literally write &getRef() if you wanted to, and it would give you the address of x.)

Rvalues are--well, rvalues are not lvalues. An expression is an rvalue if it results in a temporary object. For example:

```
1 | int x;  
2 | int getVal ()  
3 | {  
4 |     return x;  
5 | }  
6 | getVal();
```

Here, getVal() is an rvalue--the value being returned is not a reference to x, it's just a temporary value. This gets a little bit more interesting if we use real objects instead of numbers:

```
1 | string getName ()  
2 | {  
3 |     return "Alex";  
4 | }  
5 | getName();
```

Here, getName returns a string that is constructed inside the function. You can assign the result of getName to a variable:

```
1 | string name = getName();
```

But you're assigning from a temporary object, not from some value that has a fixed location. getName() is an rvalue.

## Detecting temporary objects with rvalue references

The important thing is that rvalues refer to temporary objects--just like the value returned from doubleValues. Wouldn't it be great if we could know, without a shadow of a doubt, that a value returned from an expression was temporary, and somehow write code that is overloaded to behave differently for temporary objects? Why, yes, yes indeed it would be. And this is what rvalue references are for. An rvalue reference is a reference that will bind only to a temporary object. What do I mean?

Prior to C++11, if you had a temporary object, you could use a "regular" or "lvalue reference" to bind it, but only if it was **const**:

```
1 | const string& name = getName(); // ok
2 | string& name = getName(); // NOT ok
```

The intuition here is that you cannot use a "mutable" reference because, if you did, you'd be able to modify some object that is about to disappear, and that would be dangerous. Notice, by the way, that holding on to a **const** reference to a temporary object ensures that the temporary object isn't immediately destructed. This is a nice guarantee of C++, but it is still a temporary object, so you don't want to modify it.

In C++11, however, there's a new kind of reference, an "rvalue reference", that will let you bind a mutable reference to an rvalue, but not an lvalue. In other words, rvalue references are perfect for detecting if a value is temporary object or not. Rvalue references use the **&&** syntax instead of just **&**, and can be **const** and non-**const**, just like lvalue references, although you'll rarely see a **const** rvalue reference (as we'll see, mutable references are kind of the point):

```
1 | const string&& name = getName(); // ok
2 | string&& name = getName(); // also ok -
```

So far this is all well and good, but how does it help? The most important thing about lvalue references vs rvalue references is what happens when you write functions that take lvalue or rvalue references as arguments. Let's say we have two functions:

```
1 | printReference (const String& str)
2 | {
3 |     cout << str;
4 | }
5 |
6 | printReference (String&& str)
7 | {
8 |     cout << str;
9 | }
```

Now the behavior gets interesting--the `printReference` function taking a **const** lvalue reference will accept any argument that it's given, whether it be an lvalue or an rvalue, and regardless of whether the lvalue or rvalue is mutable or not. However, in the presence of the second overload, `printReference` taking an rvalue reference, it will be given all values *except* mutable rvalue-references. In other words, if you write:

```
1 | string me( "alex" );
2 | printReference( me ); // calls the first
3 |
4 | printReference( getName() ); // calls the second
```

Now we have a way to determine if a reference variable refers to a temporary object or to a permanent object. The rvalue reference version of the method is like the secret back door entrance to the club that you can only get into if you're a temporary object (boring club, I guess). Now that we have our method of determining if an object was a temporary or a permanent thing, how can we use it?

## Move constructor and move assignment operator

The most common pattern you'll see when working with rvalue references is to create a move constructor and move assignment operator (which follows the same principles). A move constructor, like a copy constructor, takes an instance of an object as its argument and creates a new instance based on the original object. However, the move constructor can avoid memory reallocation because we know it has been provided a temporary object, so rather than copy the fields of the object, we will move them.

What does it mean to move a field of the object? If the field is a primitive type, like `int`, we just copy it. It gets more interesting if the field is a **pointer**: here, rather than allocate and initialize new memory, we can simply steal the pointer and null out the pointer in the temporary object! We know the temporary object will no longer be needed, so we can take its pointer out from under it.

Imagine that we have a simple `ArrayWrapper` class, like this:

```
1 | class ArrayWrapper
```

```

2  {
3      public:
4          ArrayWrapper (int n)
5              : _p_vals( new int[ n ] )
6              , _size( n )
7          {}
8          // copy constructor
9          ArrayWrapper (const ArrayWrapper& other)
10             : _p_vals( new int[ other._size ] )
11             , _size( other._size )
12         {
13             for ( int i = 0; i < _size;
14                 {
15                     _p_vals[ i ] = other._p_vals[ i ];
16                 }
17         }
18         ~ArrayWrapper ()
19         {
20             delete [] _p_vals;
21         }
22     private:
23     int *_p_vals;
24     int _size;
25 };

```

Notice that the copy constructor has to both allocate memory and copy every value from the array, one at a time! That's a lot of work for a copy. Let's add a move constructor and gain some massive efficiency.

```

1  class ArrayWrapper
2  {
3      public:
4          // default constructor produces a 64 element array
5          ArrayWrapper ()
6              : _p_vals( new int[ 64 ] )
7              , _size( 64 )
8          {}
9
10         ArrayWrapper (int n)
11             : _p_vals( new int[ n ] )
12             , _size( n )
13         {}
14
15         // move constructor
16         ArrayWrapper (ArrayWrapper&& other)
17             : _p_vals( other._p_vals )
18             , _size( other._size )
19         {
20             other._p_vals = NULL;
21             other._size = 0;
22         }
23
24         // copy constructor
25         ArrayWrapper (const ArrayWrapper& c)
26             : _p_vals( new int[ other._size ] )
27             , _size( other._size )
28         {
29             for ( int i = 0; i < _size; ++i
30                 {
31                     _p_vals[ i ] = other._p_vals[ i ];
32                 }
33         }
34         ~ArrayWrapper ()
35         {
36             delete [] _p_vals;
37         }
38
39     private:
40     int *_p_vals;

```

```

41 |         int _size;
42 |     };

```

Wow, the move constructor is actually simpler than the copy constructor! That's quite a feat. The main things to notice are:

1. The parameter is a non-const rvalue reference
2. other.\_p\_vals is set to NULL

The second observation explains the first--we couldn't set other.\_p\_vals to NULL if we'd taken a const rvalue reference. But why do we need to set other.\_p\_vals = NULL? The reason is the destructor--when the temporary object goes out of scope, just like all other C++ objects, its destructor will run. When its destructor runs, it will free \_p\_vals. The same \_p\_vals that we just copied! If we don't set other.\_p\_vals to NULL, the move would not really be a move--it would just be a copy that introduces a crash later on once we start using freed memory. This is the whole point of a move constructor: to avoid a copy by changing the original, temporary object!

Again, the overload rules work such that the move constructor is called only for a temporary object--and only a temporary object that can be modified. One thing this means is that if you have a function that returns a const object, it will cause the copy constructor to run instead of the move constructor--so don't write code like this:

```

1 | const ArrayWrapper getArrayWrapper (); /

```

There's still one more situation we haven't discussed how to handle in a move constructor--when we have a field that is an object. For example, imagine that instead of having a size field, we had a metadata field that looked like this:

```

1 | class Metadata
2 | {
3 | public:
4 |     Metadata (int size, const std::stri
5 |         : _name( name )
6 |         , _size( size )
7 |     {}
8 |
9 |     // copy constructor
10 |    Metadata (const Metadata& other)
11 |        : _name( other._name )
12 |        , _size( other._size )
13 |    {}
14 |
15 |    // move constructor
16 |    Metadata (Metadata&& other)
17 |        : _name( other._name )
18 |        , _size( other._size )
19 |    {}
20 |
21 |    std::string getName () const { retu
22 |    int getSize () const { return _size
23 | private:
24 |    std::string _name;
25 |    int _size;
26 | };

```

Now our array can have a name and a size, so we might have to change the definition of ArrayWrapper like so:

```

1 | class ArrayWrapper
2 | {
3 | public:
4 |     // default constructor produces a n
5 |     ArrayWrapper ()
6 |         : _p_vals( new int[ 64 ] )
7 |         , _metadata( 64, "ArrayWrapper"
8 |     {}
9 |
10 |    ArrayWrapper (int n)

```

```

11         : _p_vals( new int[ n ] )
12         , _metadata( n, "ArrayWrapper"
13     {}
14
15     // move constructor
16     ArrayWrapper (ArrayWrapper&& other)
17         : _p_vals( other._p_vals )
18         , _metadata( other._metadata )
19     {
20         other._p_vals = NULL;
21     }
22
23     // copy constructor
24     ArrayWrapper (const ArrayWrapper& c
25         : _p_vals( new int[ other._meta
26         , _metadata( other._metadata )
27     {
28         for ( int i = 0; i < _metadata.
29         {
30             _p_vals[ i ] = other._p_val
31         }
32     }
33     ~ArrayWrapper ( )
34     {
35         delete [] _p_vals;
36     }
37 private:
38     int *_p_vals;
39     MetaData _metadata;
40 };

```

Does this work? It seems very natural, doesn't it, to just call the MetaData move constructor from within the move constructor for ArrayWrapper? The problem is that this just doesn't work. The reason is simple: the value of other in the move constructor--it's an rvalue reference. But an rvalue reference is not, in fact, an rvalue. It's an lvalue, and so the copy constructor is called, not the move constructor. This is weird. I know--it's confusing. Here's the way to think about it. A rvalue is an expression that creates an object that is about to evaporate into thin air. It's on its last legs in life--or about to fulfill its life purpose. Suddenly we pass the temporary to a move constructor, and it takes on new life in the new scope. In the context where the rvalue expression was evaluated, the temporary object really is over and done with. But in our constructor, the object has a name; it will be alive for the entire duration of our function. In other words, we might use the variable other more than once in the function, and the temporary object has a defined location that truly persists for the entire function. It's an lvalue in the true sense of the term locator value, we can locate the object at a particular address that is stable for the entire duration of the function call. We might, in fact, want to use it later in the function. If a move constructor were called whenever we held an object in an rvalue reference, we might use a moved object, by accident!

```

1 // move constructor
2 ArrayWrapper (ArrayWrapper&& other)
3     : _p_vals( other._p_vals )
4     , _metadata( other._metadata )
5 {
6     // if _metadata( other._metadata ) c
7     // other._metadata here would be ext
8     other._p_vals = NULL;
9 }

```

Put a final way: both lvalue and rvalue references are lvalue expressions. The difference is that an lvalue reference must be const to hold a reference to an rvalue, whereas an rvalue reference can always hold a reference to an rvalue. It's like the difference between a pointer, and what is pointed to. The thing pointed-to came from an rvalue, but when we use rvalue reference itself, it results in an lvalue.

## std::move

So what's the trick to handling this case? We need to use std::move, from <utility>--std::move is a way of saying, "ok, honest to God I know I have an lvalue, but I want it to be an rvalue." std::move does not, in and of itself,

move anything; it just turns an lvalue into an rvalue, so that you can invoke the move constructor. Our code should look like this:

```
1  #include <utility> // for std::move
2
3  // move constructor
4  ArrayWrapper (ArrayWrapper&& other)
5      : _p_vals( other._p_vals )
6      , _metadata( std::move( other._m
7  {
8      other._p_vals = NULL;
9  }
```

And of course we should really go back to MetaData and fix its own move constructor so that it uses std::move on the string it holds:

```
1  MetaData (MetaData&& other)
2      : _name( std::move( other._name ) )
3      : _size( other._size )
4  {}
```

## Move assignment operator

Just as we have a move constructor, we should also have a move assignment operator. You can easily write one using the same techniques as for creating a move constructor.

## Move constructors and implicitly generated constructors

As you know, in C++ when you declare any constructor, the compiler will no longer generate the default constructor for you. The same is true here: adding a move constructor to a class will require you to declare and define your own default constructor. On the other hand, declaring a move constructor does not prevent the compiler from providing an implicitly generated copy constructor, and declaring a move assignment operator does not inhibit the creation of a standard assignment operator.

## How does std::move work

You might be wondering, how does one write a function like std::move? How do you get this magical property of transforming an lvalue into an rvalue reference? The answer, as you might guess, is **typecasting**. The actual declaration for std::move is somewhat more involved, but at its heart, it's just a **static\_cast** to an rvalue reference. This means, actually, that you don't really *need* to use move--but you should, since it's much more clear what you mean. The fact that a cast is required is, by the way, a very good thing! It means that you cannot accidentally convert an lvalue into an rvalue, which would be dangerous since it might allow an accidental move to take place. You must explicitly use std::move (or a cast) to convert an lvalue into an rvalue reference, and an rvalue reference will never bind to an lvalue on its own.

## Returning an explicit rvalue-reference from a function

Are there ever times where you should write a function that returns an rvalue reference? What does it mean to return an rvalue reference anyway? Aren't functions that return objects by value already rvalues?

Let's answer the second question first: returning an explicit rvalue reference is different than returning an object by value. Take the following simple example:

```
1  int x;
2
3  int getInt ()
4  {
5      return x;
6  }
7
8  int && getRvalueInt ()
9  {
10     // notice that it's fine to move a
```

```

11 |         return std::move( x );
12 |     }

```

Clearly in the first case, despite the fact that `getInt()` is an rvalue, there is a copy of the variable `x` being made. We can even see this by writing a little helper function:

```

1 | void printAddress (const int& v) // cons
2 | {
3 |     cout << reinterpret_cast<const void*>
4 | }
5 |
6 | printAddress( getInt() );
7 | printAddress( x );

```

When you run this program, you'll see that there are two separate values printed.

On the other hand,

```

1 | printAddress( getRvalueInt() );
2 | printAddress( x );

```

prints the same value because we are explicitly returning an rvalue here.

So returning an rvalue reference is a different thing than not returning an rvalue reference, but this difference manifests itself most noticeably if you have a pre-existing object you are returning instead of a temporary object created in the function (where the compiler is likely to eliminate the copy for you).

Now on to the question of whether you want to do this. The answer is: probably not. In most cases, it just makes it more likely that you'll end up with a dangling reference (a case where the reference exists, but the temporary object that it refers to has been destroyed). The issue is quite similar to the danger of returning an lvalue reference--the referred-to object may no longer exist. Rvalue references cannot magically keep an object alive for you. Returning an rvalue reference would primarily make sense in very rare cases where you have a member function and need to return the result of calling `std::move` on a field of the class from that function--and how often are you going to do that?

## Move semantics and the standard library

Going back to our original example--we were using a vector, and we don't have control over the vector class and whether or not it has a move constructor or move assignment operator. Fortunately, the standards committee is wise, and move semantics has been added to the standard library. This means that you can now efficiently return vectors, maps, strings and whatever other standard library objects you want, taking full advantage of move semantics.

### Moveable objects in STL containers

In fact, the standard library goes one step further. If you enable move semantics in your own objects by creating move assignment operators and move constructors, when you store those objects in a container, the **STL** will automatically use `std::move`, automatically taking advantage of move-enabled classes to eliminate inefficient copies.

## Move semantics and rvalue reference compiler support

Rvalue references are supported by GCC, the Intel compiler and MSVC.