

# [The C++ scientist](#)

## Scientific computing, numerical methods and optimization in C++

- [RSS](#)

<input type="text" value="Search"/>
<input type="text" value="Navigate..."/>

- [Blog](#)
- [Archives](#)
- [About](#)

## Aligned Memory Allocator

Dec 6th, 2014

### Introduction

In a previous [article about SIMD wrappers](#), I suggested to design a dedicated memory allocator to handle SIMD memory alignment constraints, but I didn't give any details on how to do it. That's the purpose of this article. The C++ standard describes a set of requirements our allocator must respect to work with standard containers. After a survey of these standard requirements, we'll see how to implement an aligned memory allocator that meets them.

### 1. The C++ standard requirements

The standard requires the allocator to define the following type:

- **value\_type**: the type of allocated elements, T
- **pointer** and **const\_pointer**: pointer and constant pointer to T
- **reference** and **constant\_reference**: reference and constant reference to T
- **size\_type**: an unsigned integral type that can represent the size of the largest object in the allocation model
- **difference\_type**: a signed integral type that can represent the difference between any two pointers in the allocation model

The standard then requires a template class **rebind** member, which is a template typedef: `Allocator<T>::rebind<U>::other` is the same type as `Allocator<U>`. This member is used by container that allocates memory for internal structures that hold T elements instead of allocating memory for T elements. For instance, `std::list<T>` allocates memory for `Node<T>` instead of T, but you don't want to use `Allocator<Node<T>>` as template parameter since this would expose implementation details in the interface. Thus you use `Allocator<T>` and the internal allocation is done with `Allocator<T>::rebind<U>::other.allocate(n)`.

Then we have to provide the **address** functions, which return the address of a given object. Two overloads are provided, one for references and one for constant references.

The two following functions are the essential part of the allocator: **allocate** and **deallocate**, which allocates/deallocates memory for n objects of type T. These functions are low-level memory management functions and are not responsible for constructing or destroying objects, this has to be done in specific functions: **construct** and **destroy**.

The last specific function is **max\_size**, a function that returns the maximum value that can be passed to allocate.

Finally, the allocator must provide default and copy constructors, and equality check operators.

### 2. Aligned memory allocator interface

Since we must handle different memory alignment bounds, our aligned memory allocator will take two template parameters: T, the type of allocated objects, and N, the alignment bound. Given the requirements of the previous section, the allocator interface looks like:

aligned\_allocator.hpp

```
1  template <class T, int N>
2      class aligned_allocator
3      {
4      {
5      public:
6
7          typedef T value_type;
8          typedef T& reference;
9          typedef const T& const_reference;
10         typedef T* pointer;
11         typedef const T* const_pointer;
12         typedef size_t size_type;
13         typedef ptrdiff_t difference_type;
14
15         template <class U>
16             struct rebind
17             {
18                 typedef aligned_allocator<U,N> other;
19             };
20
21         inline aligned_allocator() throw() {}
22         inline aligned_allocator(const aligned_allocator&) throw() {}
23
24         template <class U>
25             inline aligned_allocator(const aligned_allocator<U,N>&) throw() {}
26
27         inline ~aligned_allocator() throw() {}
28
29         inline pointer address(reference r) { return &r; }
30         inline const_pointer address(const_reference r) const { return &r; }
31
32         pointer allocate(size_type n, typename std::allocator<void>::const_pointer hint = 0);
33         inline void deallocate(pointer p, size_type);
34
35         inline void construct(pointer p, const_reference value) { new (p) value_type(value); }
36         inline void destroy(pointer p) { p->~value_type(); }
37
38         inline size_type max_size() const throw() { return size_type(-1) / sizeof(T); }
39
40         inline bool operator==(const aligned_allocator&) { return true; }
41         inline bool operator!=(const aligned_allocator& rhs) { return !operator==(rhs); }
42     };
43
```

Nothing special to say here, the **construct** function calls the copy constructor of T through the placement new operator but does not allocate memory for the element, this is the responsibility of the **allocate** function. Same thing for the **destroy** function, it calls the destructor of T but it doesn't deallocate memory, this has to be done after with a call to the **deallocate** function.

### 3. Allocate and deallocate implementation

We can now focus on the **allocate** and **deallocate** implementation. Depending on our platform, aligned memory allocation function may be available:

- On Windows 64 bits, FreeBSD (except on ARM and MIPS architectures), and Apple, the **malloc** function is already 16-bytes aligned
- Systems implementing POSIX provide the **posix\_memalign** function
- SSE intrinsics provide the **\_mm\_malloc** function
- Visual Studio provides the **\_aligned\_malloc** function

Except for the 16-bytes aligned malloc, every function takes an alignment parameter that must be a power of 2, thus the N template parameter of our allocator should be a power of 2 so it can work with these aligned memory allocation functions. Note that many of these functions can be available on a same platform.

Assume we can detect at compile time if such functions are available (we'll come back on this later); we can provide a function that selects the aligned memory allocation to use if such a function is available, otherwise forwards to our own implementation:

aligned\_allocator.hpp

```

1 namespace detail
2 {
3     inline void* _aligned_malloc(size_t size, size_t alignment)
4     {
5         void* res = 0;
6         void* ptr = malloc(size+alignment);
7         if(ptr != 0)
8         {
9             res = reinterpret_cast<void*>((reinterpret_cast<size_t>(ptr) & ~(size_t)(alignment-1))) + alignment);
10            *(reinterpret_cast<void**>(res) - 1) = ptr;
11        }
12        return res;
13    }
14 }
15
16 inline void* aligned_malloc(size_t size, size_t alignment)
17 {
18     #if MALLOC_ALREADY_ALIGNED
19         return malloc(size);
20     #elif HAS_MM_MALLOC
21         return _mm_malloc(size,alignment);
22     #elif HAS_POSIX_MEMALIGN
23         void* res;
24         const int failed = posix_memalign(&res,size,alignment);
25         if(failed) res = 0;
26         return res;
27     #elif (defined _MSC_VER)
28         return _aligned_malloc(size, alignment);
29     #else
30         return detail::_aligned_malloc(size,alignment);
31     #endif
32 }
33

```

The idea in the **\_aligned\_malloc** function is to search for the first aligned memory address (*res*) after the one returned by the classic malloc function (*ptr*), and to use it as return value. But since we must ensure *size* bytes are available after *res*, we must allocate more than *size* bytes; the minimum size to allocate to prevent buffer overflow is *size+alignment*. Then we store the *ptr* value just before *res* so the **\_aligned\_free** function can easily retrieve and pass it to the classic free function:

aligned\_allocator.hpp

```

1 namespace detail
2 {
3     inline void _aligned_free(void* ptr)
4     {
5         if(ptr != 0)
6             free(*(reinterpret_cast<void**>(ptr)-1));
7     }
8 }
9
10 inline void aligned_free(void* ptr)
11 {
12     #if MALLOC_ALREADY_ALIGNED
13         free(ptr);
14     #elif HAS_MM_MALLOC
15         _mm_free(ptr);
16     #elif HAS_POSIX_MEMALIGN
17         free(ptr);
18     #elif defined(_MSC_VER)
19         _aligned_free(ptr);
20     #else
21         detail::_aligned_free(ptr);
22     #endif
23 }
24

```

The **aligned\_free** function is the symmetric of **aligned\_malloc**: it selects the aligned memory function available or forwards to the **\_aligned\_free** function.

We can now write the **allocate** and **deallocate** functions of the allocator:

aligned\_allocator.hpp

```
1 template <class T, int N>
2     typename aligned_allocator<T,N>::pointer
3     aligned_allocator<T,N>::allocate(size_type n, typename std::allocator<void>::const_pointer hint)
4     {
5         pointer res = reinterpret_cast<pointer>(aligned_malloc(sizeof(T)*n,N));
6         if(res == 0)
7             throw std::bad_alloc();
8         return res;
9     }
10 }
11 template <class T, int N>
12     typename aligned_allocator<T,N>::pointer
13     aligned_allocator<T,N>::deallocate(pointer p, size_type)
14     {
15         aligned_free(p);
16     }
17 }
```

Here we see the advantage to have encapsulated aligned memory allocation selection in a dedicated function: the allocate function of the allocator simply forwards to this dedicated function and then handles possible bad allocation. The result is a simple and easy-to-read code. Another advantage is that you can use **aligned\_malloc** and **aligned\_free** functions outside the **aligned\_allocator** class if you need.

Note: the call to malloc after the MALLOC\_ALREADY\_ALIGNED preprocessor token should be available for 16-bytes aligned memory allocator only (the same applies to the call to free). Thus we should provide two versions of **aligned\_malloc** and **aligned\_free** and a specialization of the allocator off **N = 16**.

## 4. Detecting available aligned memory allocation

Now that we have implemented the allocation and deallocation methods, we can come back to the preprocessor tokens. Defining these tokens is not simple because you have to refer to the documentation of a lot of various systems and architectures. Thus there's a chance that we may not be comprehensive, but at least we can cover the most common platforms.

Let's start with the GNU world; according to this [documentation](#), "The address of a block returned by malloc or realloc in GNU systems is always a multiple of eight (or sixteen on 64-bit systems)". According to this [one](#), page 114, "[The] LP64 model [...] is used by all 64-bit UNIX ports", therefore we should use this predefined macro instead of `__x86_64__` (this last one won't work on PowerPC or SPARC). Thus we can define the following macro:

aligned\_allocator.hpp

```
1 #if defined(__GLIBC__) && ((__GLIBC__>=2 && __GLIBC_MINOR__ >= 8) || __GLIBC__>2) \
2     && defined(__LP64__)
3     #define GLIBC_MALLOC_ALREADY_ALIGNED 1
4 #else
5     #define GLIBC_MALLOC_ALREADY_ALIGNED 0
6 #endif
7
```

FreeBSD has 16-byte aligned malloc, except on ARM and MIPS architectures (see this [documentation](#)):

aligned\_allocator.hpp

```
1 #if defined(__FreeBSD__) && !defined(__arm__) && !defined(__mips__)
2     #define FREEBSD_MALLOC_ALREADY_ALIGNED 1
3 #else
4     #define FREEBSD_MALLOC_ALREADY_ALIGNED 0
5 #endif
6
```

On windows 64 bits and Apple OS, the malloc function is also already aligned, thus we can define the **MALLOC\_ALREADY\_ALIGNED** macro, based on these information and the macros previously defined:

aligned\_allocator.hpp

```
1 #if (defined(__APPLE__) \
2     || defined(_WIN64) \
```

```

3 || GLIBC_MALLOC_ALREADY_ALIGNED \
4 || FreeBSD_MALLOC_ALREADY_ALIGNED)
5 #define MALLOC_ALREADY_ALIGNED 1
6 #else
7 #define MALLOC_ALREADY_ALIGNED 0
8 #endif
9

```

To handle systems implementing POSIX:

aligned\_allocator.hpp

```

1 #if ((defined __QNXNTO__) || (defined _GNU_SOURCE) || ((defined _XOPEN_SOURCE) && (_XOPEN_SOURCE >= 600))) \
2 && (defined _POSIX_ADVISORY_INFO) && (_POSIX_ADVISORY_INFO > 0)
3 #define HAS_POSIX_MEMALIGN 1
4 #else
5 #define HAS_POSIX_MEMALIGN 0
6 #endif
7

```

The last macro to define is **HAS\_MM\_MALLOC**; the `_mm_malloc` function is provided with SSE intrinsics, thus we can rely on the macros defined in this [article](#):

aligned\_allocator.hpp

```

1 #if SSE_INSTR_SET > 0
2 #define HAS_MM_MALLOC 1
3 #else
4 #define HAS_MM_MALLOC 0
5 #endif
6

```

That's it, some architectures may be missing but it shouldn't be too complicated to handle them with appropriate documentation.

## Conclusion

The aligned memory allocator designed in this article meets the standard requirements and can therefore be used with any container of the STL. If you work with intrinsics wrappers and `std::vector`, this allocator will allow you to load the data from memory with the [load\\_a](#) method, faster than [load\\_u](#) (the same applies for storing data to memory):

sample.cpp

```

1 typedef std::vector<double,aligned_allocator<double,16> > vector_type;
2 vector_type v1,v2,v3;
3 // code filling v1 and v2
4 for(size_t i = 0; i < v1.size(); i += simd_traits<double>::size)
5 {
6     vector2d v1d = load_a(&v1[i]);
7     vector2d v2d = load_a(&v2[i]);
8     vector2d v3d = v1d + v2d;
9     store_a(&v3[i],v3d);
10 }
11

```

But as we will see in a forthcoming article, `std::vector` may not be the most appropriate container for efficient numerical analysis.

Posted by Johan Mabilie Dec 6th, 2014 [Memory](#)

[« Performance considerations about SIMD wrappers](#)

## Comments