

[The C++ scientist](#)

Scientific computing, numerical methods and optimization in C++

- [RSS](#)

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

Writing C++ Wrappers for SIMD Intrinsics (5)

Oct 25th, 2014

4. Making the code more generic

In the previous section we saw how to plug the wrappers into existing code and ended up with the following loop:

sample.cpp

```
1 std::vector<float> a, b, c, d, e
2 // Somewhere in the code the vectors are resized
3 // so they hold n elements
4 for(size_t i = 0; i < n/4; i+=4)
5 {
6     vector4f av; av.load_a(&a[i]);
7     vector4f bv; bv.load_a(&b[i]);
8     vector4f cv; cv.load_a(&c[i]);
9     vector4f dv; dv.load_a(&d[i]);
10
11     vector4f ev = av*bv + cv*dv;
12     ev.store_a(&e[i]);
13 }
14 // Remaining part of the loop
15 // ...
16
```

As said in the previous section, the first problem of this code is its lack of genericity; we are highly coupled with the SIMD instruction set wrapped, and replacing it with another one requires code changes we should avoid. If we want to make the code independant from the SIMD instruction set and the related wrapper, we need to hide the specifics of this instruction set, that is, the vector type and its size (the number of scalars it holds).

4.1 Hiding the wrapper type

We want to be able to select the right wrapper depending on the scalar type and the instruction set used. When talking about selecting a type depending on another one, the first thing that comes to mind is type traits. Here our traits must contain the wrapper type and its size associated with the scalar type used:

simd_traits.hpp

```
1 template <class T>
2 struct simd_traits
```

```

3 {
4     typedef T type;
5     static const size_t size = 1;
6 };
7

```

The general definition of the traits class allows us to write code that works even for types that don't have related wrappers (numerical types defined by another user for instance). Then we need to specialize these definitions for float and double, depending on the considered instruction set. Assume we can detect the instruction set available on our system and save this information in a macro (we'll see how to do that in a later section). The specialization of the traits class will look like:

simd.hpp

```

1  #ifndef USE_SSE
2  template <>
3      struct simd_traits<float>
4      {
5          typedef vector4f type;
6          static const size_t size = 4;
7      };
8
9  template <>
10     struct simd_traits<double>
11     {
12         typedef vector2d type;
13         static const size_t size = 2;
14     };
15 #elif USE_AVX
16 template <>
17     struct simd_traits<float>
18     {
19         typedef vector8f type;
20         static const size_t size = 8;
21     };
22
23 template <>
24     struct simd_traits<double>
25     {
26         typedef vector4d type;
27         static const size_t size = 4;
28     };
29 #endif
30

```

Now we can adapt the loop so it doesn't explicitly refer to the vector4f type:

sample.cpp

```

1  std::vector<float> a,b,c,d,e;
2  // ... resize a, b, c, d, and e so they hold n elements
3  typedef simd_traits<float>::type vec_type;
4  size_t vec_size = simd_traits<float>::size;
5  for(size_t i = 0; i < n/vec_size; i += vec_size)
6  {
7      vec_type av; av.load_a(&a[i]);
8      vec_type bv; bv.load_a(&b[i]);
9      vec_type cv; cv.load_a(&c[i]);
10     vec_type dv; dv.load_a(&d[i]);
11
12     vec_type ev = av*bv + cv*dv;
13     ev.store_a(&e[i]);
14 }
15 // Remaining part of the loop
16 // ...
17

```

That's it! If we need to compile this code on a system where AVX is available, we have nothing to do. The macro `USE_AVX` will be defined, the specialization of `simd_traits` with `vector8f` as inner type will be instantiated, and the loop will use the `vector8f` wrapper and the AVX intrinsics. However, there's still a problem: we can migrate to any SIMD instruction set for which a wrapper is available, but we can't use types that don't have related wrappers. The `simd_traits` works fine even for user defined types, but the load and store functions are available for wrappers only. We need to provide generic versions of these functions that work with any type.

4.2 Generic load and store functions

Actually, all we have to do is to provide two versions of these functions: one for types that don't have related wrappers, and one that works with wrappers. Template specialization can be of help here, but since partial specialization is not possible for functions, let's wrap them into a `simd_functions_invoker` class:

`simd.hpp`

```
1 // Common implementation for types that support vectorization
2 template <class T, class V>
3     struct simd_functions_invoker
4     {
5         inline static V
6         set1(const T& a) { return V(a); }
7
8         inline static V
9         load_a(const T* src) { V res; res.load_a(src); return res; }
10
11         inline static V
12         load_u(const T* src) { V res; res.load_u(src); return res; }
13
14         inline static void
15         store_a(T* dst, const V& src) { src.store_a(dst); }
16
17         inline static void
18         store_u(T* dst, const V& src) { src.store_u(dst); }
19     };
20
21 // Specialization for types that don't support vectorization
22 template <class T>
23     struct simd_functions_invoker<T,T>
24     {
25         inline static T
26         set1(const T& a) { return T(a); }
27
28         inline static T
29         load_a(const T* src) { return *src; }
30
31         inline static T
32         load_u(const T* src) { return *src; }
33
34         inline static void
35         store_a(T* dst, const T& src) { *dst = src; }
36
37         inline static void
38         store_u(T* dst, const T& src) { *dst = src; }
39     };
40
```

We've added the `set1` function so we can initialize wrappers and scalar type from a single value in an uniform way. Calling the generic functions would look like:

`sample.cpp`

```
1 typedef simd_traits<float>::simd_type vec_type;
2 vec_type va = simd_functions_invoker<float,vec_type>::load_a(a);
3
```

That's too much verbose. Let's add façade functions that deduce template parameters for us:

simd.hpp

```
1 template <class T> inline typename simd_traits<T>::type
2 set1(const T& a)
3 { return simd_functions_invoker<T,typename simd_traits<T>::type>::set1(a); }
4
5 template <class T> inline typename simd_traits<T>::type
6 load_a(const T* src)
7 { return simd_functions_invoker<T,typename simd_traits<T>::type>::load_a(src); }
8
9 template <class T> inline typename simd_traits<T>::type
10 load_u(const T* src)
11 { return simd_functions_invoker<T,typename simd_traits<T>::type>::load_u(src); }
12
13 template <class T> inline void
14 store_a(T* dst, const typename simd_traits<T>::type& src)
15 { simd_functions_invoker<T,typename simd_traits<T>::type>::store_a(dst,src); }
16
17 template <class T> inline void
18 store_u(T* dst, const typename simd_traits<T>::type& src)
19 { simd_functions_invoker<T,typename simd_traits<T>::type>::store_u(dst,src); }
20
```

Now we can use these generic functions in the previous loop so it works with any type, even those that don't support vectorization:

sample.cpp

```
1 std::vector<float> a,b,c,d,e;
2 // ... resize a, b, c, d, and e so they hold n elements
3 typedef simd_traits<float>::type vec_type;
4 size_t vec_size = simd_traits<float>::size;
5 for(size_t i = 0; i < n/vec_size; i += vec_size)
6 {
7     vec_type av = load_a(&a[i]);
8     vec_type bv = load_a(&b[i]);
9     vec_type cv = load_a(&c[i]);
10    vec_type dv = load_a(&d[i]);
11
12    vec_type ev = av*bv + cv*dv;
13    store_a(&e[i],ev);
14 }
15 // Remaining part of the loop
16 // ...
17
```

Or, if you want to be more concise:

sample.cpp

```
1 std::vector<float> a,b,c,d,e;
2 // ... resize a, b, c, d, and e so they hold n elements
3 typedef simd_traits<float>::type vec_type;
4 size_t vec_size = simd_traits<float>::size;
5 for(size_t i = 0; i < n/vec_size; i += vec_size)
6 {
7     vec_type ev = load_a(&a[i])*load_a(&b[i]) + load_a(&c[i])*load_a(&d[i]));
8     store_a(&e[i], ev);
9 }
10 // Remaining part of the loop
11 // ...
12
```

We've reached our goal, we can use intrinsics almost like floats; in a real application code, it is likely that you initialize the wrappers through load functions, then perform the computations and finally store the results (like in the not concise version of the generic loop); thus the only difference between classical code and code with SIMD wrappers is the initialization and storing of wrappers (and eventually the functions signatures if you want to pass wrappers instead of scalars), the other parts should be exactly the same and the code remains easy to read and to maintain.

4.3 Detecting the supported instruction set

Until now, we've assumed we were able to detect at compile time the available instruction sets. Let's see now how to achieve this. Compilers often provide preprocessor tokens depending on the available instruction sets, but these tokens may vary from one compiler to another, so we have to standardize that. On most 64-bit compilers, the tokens look like `__SSE__` or `__SSE3__`, on 32-bit systems, Microsoft compilers set the preprocessor token `_M_IX86_FP` to 1 for SSE (vectorization of float) and 2 for SSE2 (vectorization of double and integers).

Here is how we can standardize that:

`simd_config.hpp`

```
1  #if (defined(_M_AMD64) || defined(_M_X64) || defined(__amd64)) && ! defined(__x86_64__)
2      #define __x86_64__ 1
3  #endif
4
5  // Find sse instruction set from compiler macros if SSE_INSTR_SET not defined
6  // Note: Not all compilers define these macros automatically
7  #ifndef SSE_INSTR_SET
8      #if defined ( __AVX2__ )
9          #define SSE_INSTR_SET 8
10         #elif defined ( __AVX__ )
11             #define SSE_INSTR_SET 7
12         #elif defined ( __SSE4_2__ )
13             #define SSE_INSTR_SET 6
14         #elif defined ( __SSE4_1__ )
15             #define SSE_INSTR_SET 5
16         #elif defined ( __SSSE3__ )
17             #define SSE_INSTR_SET 4
18         #elif defined ( __SSE3__ )
19             #define SSE_INSTR_SET 3
20         #elif defined ( __SSE2__ ) || defined ( __x86_64__ )
21             #define SSE_INSTR_SET 2
22         #elif defined ( __SSE__ )
23             #define SSE_INSTR_SET 1
24         #elif defined ( _M_IX86_FP ) // Defined in MS compiler on 32bits system. 1: SSE, 2: SSE2
25             #define SSE_INSTR_SET _M_IX86_FP
26         #else
27             #define SSE_INSTR_SET 0
28         #endif // instruction set defines
29 #endif // SSE_INSTR_SET
30
```

Now we can use the `SSE_INSTR_SET` token to include the right file:

`simd_config.hpp`

```
1  // Include the appropriate header file for intrinsic functions
2  #if SSE_INSTR_SET > 7 // AVX2 and later
3      #ifdef __GNUC__
4          #include <x86intrin.h> // x86intrin.h includes header files for whatever instruction
5                                  // sets are specified on the compiler command line, such as:
6                                  // xopintrin.h, fma4intrin.h
7      #else
8          #include <immintrin.h> // MS version of immintrin.h covers AVX, AVX2 and FMA3
9      #endif // __GNUC__
10 #elif SSE_INSTR_SET == 7
11     #include <immintrin.h> // AVX
12 #elif SSE_INSTR_SET == 6
13     #include <nmmintrin.h> // SSE4.2

```

```

14 #elif SSE_INSTR_SET == 5
15     #include <smmintrin.h>           // SSE4.1
16 #elif SSE_INSTR_SET == 4
17     #include <tmmintrin.h>          // SSE3
18 #elif SSE_INSTR_SET == 3
19     #include <pmmintrin.h>          // SSE3
20 #elif SSE_INSTR_SET == 2
21     #include <emmintrin.h>          // SSE2
22 #elif SSE_INSTR_SET == 1
23     #include <xmmmintrin.h>         // SSE
24

```

Note that if you split the implementation of SSE wrappers and AVX wrappers into different files, you can also use the `SSE_INSTR_SET` token to include the implementation file in the `simd.hpp` file:

`simd.hpp`

```

1  #include "simd_config.hpp"
2  #if SSE_INSTR_SET > 6
3      #include "simd_avx.hpp"
4  #endif
5  #if SSE_INSTR_SET > 0
6      #include "simd_sse.hpp"
7  #endif
8
9  // Definition of traits and generic load and store functions
10 // ...
11

```

Now from the client code, the only file to include is `simd.hpp`, and everything will be available.

4.4 Going further

Now that we have nice wrappers providing basic functionalities, what could be the next step ? Well, first we could add a method to retrieve an element in the vector:

`simd_base.hpp`

```

1  template <class X>
2      class simd_vector
3      {
4      public:
5
6          typedef simd_traits<X>::value_type value_type;
7
8          // ...
9
10         value_type operator[](size_t index) const
11         {
12             size_t size = simd_traits<X>::size;
13             value_type v[size];
14             (*this)().store_u(v);
15             return v[index];
16         }
17     };
18

```

We can add horizontal add function, useful for linear algebra products:

`simd_sse.hpp`

```

1  inline float hadd(const vector4f& rhs)
2  {
3      #if SSE_INSTR_SET >= 3 // SSE3
4          __m128 tmp0 = _mm_hadd_ps(rhs,rhs);

```

```
5     __m128 tmp1 = _mm_hadd_ps(tmp0,tmp0);  
6     #else  
7         __m128 tmp0 = _mm_add_ps(rhs,_mm_movehl_ps(rhs,rhs));  
8         __m128 tmp1 = _mm_add_ss(tmp0,_mm_shuffle_ps(tmp0,tmp0,1));  
9     #endif  
10    return _mm_cvtss_f32(tmp1);  
11    }  
12
```

Another useful project would be to write overloads of standard mathematical functions (exp, log, etc) that work with the wrappers.

As you can see, writing the wrappers is just the beginning, you can then enrich them with whatever functionality you need but this goes beyond the topic of this first series of articles.

Posted by Johan Mabilie Oct 25th, 2014 [SIMD](#), [Vectorization](#)

[« Writing C++ wrappers for SIMD intrinsics \(4\) Performance considerations about SIMD wrappers »](#)

Comments