

---

# Boost.Multiprecision

John Maddock

Christopher Kormanyos

Copyright © 2002-2013 John Maddock and Christopher Kormanyos

Distributed under the Boost Software License, Version 1.0. (See accompanying file LICENSE\_1\_0.txt or copy at [http://www.boost.org/LICENSE\\_1\\_0.txt](http://www.boost.org/LICENSE_1_0.txt))

## Table of Contents

Introduction .....	3
Tutorial .....	9
Integer Types .....	9
cpp_int .....	9
gmp_int .....	12
tom_int .....	13
Examples .....	15
Factorials .....	15
Bit Operations .....	16
Floating Point Numbers .....	17
cpp_bin_float .....	18
cpp_dec_float .....	20
gmp_float .....	21
mpfr_float .....	23
float128 .....	25
Examples .....	26
Area of Circle .....	26
Defining a Special Function. ....	27
Calculating a Derivative .....	30
Calculating an Integral .....	32
Polynomial Evaluation .....	34
Interval Number Types .....	36
mpfi_float .....	36
Rational Number Types .....	39
cpp_rational .....	40
gmp_rational .....	41
tommath_rational .....	42
Use With Boost.Rational .....	43
rational_adaptor .....	43
Miscellaneous Number Types. ....	44
logged_adaptor .....	44
debug_adaptor .....	46
Visual C++ Debugger Visualizers .....	48
Constructing and Interconverting Between Number Types .....	51
Generating Random Numbers .....	53
Primality Testing .....	55
Literal Types and constexpr Support .....	56
Rounding Rules for Conversions .....	58
Mixed Precision Arithmetic .....	59
Generic Integer Operations .....	61
Boost.Serialization Support .....	62
Numeric Limits .....	62

std::numeric_limits<> constants .....	63
std::numeric_limits<> functions .....	70
Numeric limits for 32-bit platform .....	75
How to Determine the Kind of a Number From std::numeric_limits .....	79
Input Output .....	81
Reference .....	84
number .....	84
cpp_int .....	97
gmp_int .....	98
tom_int .....	99
gmp_float .....	99
mpfr_float_backend .....	100
cpp_bin_float .....	100
cpp_dec_float .....	102
Internal Support Code .....	102
Backend Requirements .....	103
Header File Structure .....	126
Performance Comparison .....	128
The Overhead in the Number Class Wrapper .....	128
Floating-Point Real World Tests .....	128
Integer Real World Tests .....	129
Float Algorithm Performance .....	131
Integer Algorithm Performance .....	135
Rational Type Performance .....	143
Roadmap .....	147
History .....	147
TODO .....	149
FAQ .....	150
Acknowledgements .....	151
Indexes .....	152

# Introduction

The Multiprecision Library provides [integer](#), [rational](#) and [floating-point](#) types in C++ that have more range and precision than C++'s ordinary built-in types. The big number types in Multiprecision can be used with a wide selection of basic mathematical operations, elementary transcendental functions as well as the functions in Boost.Math. The Multiprecision types can also interoperate with the built-in types in C++ using clearly defined conversion rules. This allows Boost.Multiprecision to be used for all kinds of mathematical calculations involving integer, rational and floating-point types requiring extended range and precision.

Multiprecision consists of a generic interface to the mathematics of large numbers as well as a selection of big number back ends, with support for integer, rational and floating-point types. Boost.Multiprecision provides a selection of back ends provided off-the-rack in including interfaces to GMP, MPFR, MPIR, TomMath as well as its own collection of Boost-licensed, header-only back ends for integers, rationals and floats. In addition, user-defined back ends can be created and used with the interface of Multiprecision, provided the class implementation adheres to the necessary [concepts](#).

Depending upon the number type, precision may be arbitrarily large (limited only by available memory), fixed at compile time (for example 50 or 100 decimal digits), or a variable controlled at run-time by member functions. The types are expression-template-enabled for better performance than naive user-defined types.

The Multiprecision library comes in two distinct parts:

- An expression-template-enabled front-end `number` that handles all the operator overloading, expression evaluation optimization, and code reduction.
- A selection of back-ends that implement the actual arithmetic operations, and need conform only to the reduced interface requirements of the front-end.

Separation of front-end and back-end allows use of highly refined, but restricted license libraries where possible, but provides Boost license alternatives for users who must have a portable unconstrained license. Which is to say some back-ends rely on 3rd party libraries, but a header-only Boost license version is always available (if somewhat slower).

Should you just wish to cut to the chase and use a fully Boost-licensed number type, then skip to [cpp\\_int](#) for multiprecision integers, [cpp\\_dec\\_float](#) for multiprecision floating point types and [cpp\\_rational](#) for rational types.

The library is often used via one of the predefined typedefs: for example if you wanted an [arbitrary precision](#) integer type using [GMP](#) as the underlying implementation then you could use:

```
#include <boost/multiprecision/gmp.hpp> // Defines the wrappers around the GMP library's types

boost::multiprecision::mpz_int myint;    // Arbitrary precision integer type.
```

Alternatively, you can compose your own multiprecision type, by combining `number` with one of the predefined back-end types. For example, suppose you wanted a 300 decimal digit floating-point type based on the [MPFR](#) library. In this case, there's no predefined typedef with that level of precision, so instead we compose our own:

```
#include <boost/multiprecision/mpfr.hpp> // Defines the Backend type that wraps MPFR

namespace mp = boost::multiprecision;    // Reduce the typing a bit later...

typedef mp::number<mp::mpfr_float_backend<300> > my_float;

my_float a, b, c; // These variables have 300 decimal digits precision
```

We can repeat the above example, but with the expression templates disabled (for faster compile times, but slower runtimes) by passing a second template argument to `number`:

```
#include <boost/multiprecision/mpfr.hpp> // Defines the Backend type that wraps MPFR

namespace mp = boost::multiprecision;    // Reduce the typing a bit later...

typedef mp::number<mp::mpfr_float_backend<300>, et_off> my_float;

my_float a, b, c; // These variables have 300 decimal digits precision
```

We can also mix arithmetic operations between different types, provided there is an unambiguous implicit conversion from one type to the other:

```
#include <boost/multiprecision/cpp_int.hpp>

namespace mp = boost::multiprecision;    // Reduce the typing a bit later...

mp::int128_t a(3), b(4);
mp::int512_t c(50), d;

d = c * a; // OK, result of mixed arithmetic is an int512_t
```

Conversions are also allowed:

```
d = a; // OK, widening conversion.
d = a * b; // OK, can convert from an expression template too.
```

However conversions that are inherently lossy are either declared explicit or else forbidden altogether:

```
d = 3.14; // Error implicit conversion from float not allowed.
d = static_cast<mp::int512_t>(3.14); // OK explicit construction is allowed
```

Mixed arithmetic will fail if the conversion is either ambiguous or explicit:

```
number<cpp_int_backend<>, et_off> a(2);
number<cpp_int_backend<>, et_on> b(3);

b = a * b; // Error, implicit conversion could go either way.
b = a * 3.14; // Error, no operator overload if the conversion would be explicit.
```

## Move Semantics

On compilers that support rvalue-references, class `number` is move-enabled if the underlying backend is.

In addition the non-expression template operator overloads (see below) are move aware and have overloads that look something like:

```
template <class B>
number<B, et_off> operator + (number<B, et_off>&& a, const number<B, et_off>& b)
{
    return std::move(a += b);
}
```

These operator overloads ensure that many expressions can be evaluated without actually generating any temporaries. However, there are still many simple expressions such as:

```
a = b * c;
```

Which don't noticeably benefit from move support. Therefore, optimal performance comes from having both move-support, and expression templates enabled.

Note that while "moved-from" objects are left in a sane state, they have an unspecified value, and the only permitted operations on them are destruction or the assignment of a new value. Any other operation should be considered a programming error and all of our backends will trigger an assertion if any other operation is attempted. This behavior allows for optimal performance on move-construction (i.e. no allocation required, we just take ownership of the existing object's internal state), while maintaining usability in the standard library containers.

## Expression Templates

Class `number` is expression-template-enabled: that means that rather than having a multiplication operator that looks like this:

```
template <class Backend>
number<Backend> operator * (const number<Backend>& a, const number<Backend>& b)
{
    number<Backend> result(a);
    result *= b;
    return result;
}
```

Instead the operator looks more like this:

```
template <class Backend>
unmentionable-type operator * (const number<Backend>& a, const number<Backend>& b);
```

Where the "unmentionable" return type is an implementation detail that, rather than containing the result of the multiplication, contains instructions on how to compute the result. In effect it's just a pair of references to the arguments of the function, plus some compile-time information that stores what the operation is.

The great advantage of this method is the *elimination of temporaries*: for example the "naive" implementation of `operator*` above, requires one temporary for computing the result, and at least another one to return it. It's true that sometimes this overhead can be reduced by using move-semantics, but it can't be eliminated completely. For example, lets suppose we're evaluating a polynomial via Horner's method, something like this:

```
T a[7] = { /* some values */ };
//....
y = (((((a[6] * x + a[5]) * x + a[4]) * x + a[3]) * x + a[2]) * x + a[1]) * x + a[0]);
```

If type `T` is a number, then this expression is evaluated *without creating a single temporary value*. In contrast, if we were using the `mpfr_class` C++ wrapper for MPFR - then this expression would result in no less than 11 temporaries (this is true even though `mpfr_class` does use expression templates to reduce the number of temporaries somewhat). Had we used an even simpler wrapper around MPFR like `mpreal` things would have been even worse and no less that 24 temporaries are created for this simple expression (note - we actually measure the number of memory allocations performed rather than the number of temporaries directly, note also that the `mpf_class` wrapper that will be supplied with GMP-5.1 reduces the number of temporaries to pretty much zero). Note that if we compile with expression templates disabled and rvalue-reference support on, then actually still have no wasted memory allocations as even though temporaries are created, their contents are moved rather than copied.<sup>1</sup>

---

<sup>1</sup> The actual number generated will depend on the compiler, how well it optimises the code, and whether it supports rvalue references. The number of 11 temporaries was generated with Visual C++ 10



## Important

Expression templates can radically reorder the operations in an expression, for example:

```
a = (b * c) * a;
```

Will get transformed into:

```
a *= c; a *= b;
```

If this is likely to be an issue for a particular application, then they should be disabled.

This library also extends expression template support to standard library functions like `abs` or `sin` with number arguments. This means that an expression such as:

```
y = abs(x);
```

can be evaluated without a single temporary being calculated. Even expressions like:

```
y = sin(x);
```

get this treatment, so that variable 'y' is used as "working storage" within the implementation of `sin`, thus reducing the number of temporaries used by one. Of course, should you write:

```
x = sin(x);
```

Then we clearly can't use `x` as working storage during the calculation, so then a temporary variable is created in this case.

Given the comments above, you might be forgiven for thinking that expression-templates are some kind of universal-panacea: sadly though, all tricks like this have their downsides. For one thing, expression template libraries like this one, tend to be slower to compile than their simpler cousins, they're also harder to debug (should you actually want to step through our code!), and rely on compiler optimizations being turned on to give really good performance. Also, since the return type from expressions involving numbers is an "unmentionable implementation detail", you have to be careful to cast the result of an expression to the actual number type when passing an expression to a template function. For example, given:

```
template <class T>
void my_proc(const T&);
```

Then calling:

```
my_proc(a+b);
```

Will very likely result in obscure error messages inside the body of `my_proc` - since we've passed it an expression template type, and not a number type. Instead we probably need:

```
my_proc(my_number_type(a+b));
```

Having said that, these situations don't occur that often - or indeed not at all for non-template functions. In addition, all the functions in the Boost.Math library will automatically convert expression-template arguments to the underlying number type without you having to do anything, so:

```
mpfr_float_100 a(20), delta(0.125);
boost::math::gamma_p(a, a + delta);
```

Will work just fine, with the `a + delta` expression template argument getting converted to an `mpfr_float_100` internally by the Boost.Math library.

One other potential pitfall that's only possible in C++11: you should never store an expression template using:

```
auto my_expression = a + b - c;
```

unless you're absolutely sure that the lifetimes of `a`, `b` and `c` will outlive that of `my_expression`.

And finally... the performance improvements from an expression template library like this are often not as dramatic as the reduction in number of temporaries would suggest. For example if we compare this library with `mpfr_class` and `mpreal`, with all three using the underlying MPFR library at 50 decimal digits precision then we see the following typical results for polynomial execution:

**Table 1. Evaluation of Order 6 Polynomial.**

Library	Relative Time	Relative number of memory allocations
number	1.0 (0.00957s)	1.0 (2996 total)
<code>mpfr_class</code>	1.1 (0.0102s)	4.3 (12976 total)
<code>mpreal</code>	1.6 (0.0151s)	9.3 (27947 total)

As you can see, the execution time increases a lot more slowly than the number of memory allocations. There are a number of reasons for this:

- The cost of extended-precision multiplication and division is so great, that the times taken for these tend to swamp everything else.
- The cost of an in-place multiplication (using `operator* =`) tends to be more than an out-of-place `operator*` (typically `operator* =` has to create a temporary workspace to carry out the multiplication, where as `operator*` can use the target variable as workspace). Since the expression templates carry out their magic by converting out-of-place operators to in-place ones, we necessarily take this hit. Even so the transformation is more efficient than creating the extra temporary variable, just not by as much as one would hope.

Finally, note that `number` takes a second template argument, which, when set to `et_off` disables all the expression template machinery. The result is much faster to compile, but slower at runtime.

We'll conclude this section by providing some more performance comparisons between these three libraries, again, all are using MPFR to carry out the underlying arithmetic, and all are operating at the same precision (50 decimal digits):

**Table 2. Evaluation of Boost.Math's Bessel function test data**

Library	Relative Time	Relative Number of Memory Allocations
<code>mpfr_float_50</code>	1.0 (5.78s)	1.0 (1611963)
<code>number&lt;mpfr_float_backend&lt;50&gt;, et_off&gt;</code> (but with rvalue reference support)	1.1 (6.29s)	2.64 (4260868)
<code>mpfr_class</code>	1.1 (6.28s)	2.45 (3948316)
<code>mpreal</code>	1.65 (9.54s)	8.21 (13226029)

**Table 3. Evaluation of Boost.Math's Non-Central T distribution test data**

Library	Relative Time	Relative Number of Memory Allocations
number	1.0 (263s)	1.0 (127710873)
number<mpfr_float_backend<50>, et_off> (but with rvalue reference support)	1.0 (260s)	1.2 (156797871)
<a href="#">mpfr_class</a>	1.1 (287s)	2.1 (268336640)
<a href="#">mpreal</a>	1.5 (389s)	3.6 (466960653)

The above results were generated on Win32 compiling with Visual C++ 2010, all optimizations on (/Ox), with MPFR 3.0 and MPIR 2.3.0.



# Tutorial

In order to use this library you need to make two choices:

- What kind of number do I want ([integer](#), [floating point](#) or [rational](#)).
- Which back-end do I want to perform the actual arithmetic (Boost-supplied, GMP, MPFR, Tommath etc)?

## Integer Types

The following back-ends provide integer arithmetic:

Backend Type	Header	Radix	Dependencies	Pros	Cons
cpp_int	boost/multiprecision/cpp_int.hpp	2	None	Very versatile, Boost licensed, all C++ integer type which support both <a href="#">arbitrary precision</a> and fixed precision integer types.	Slower than <a href="#">GMP</a> , though typically not as slow as <a href="#">libtommath</a>
gmp_int	boost/multiprecision/gmp.hpp	2	<a href="#">GMP</a>	Very fast and efficient back-end.	Dependency on GNU licensed <a href="#">GMP</a> library.
tom_int	boost/multiprecision/tommath.hpp	2	<a href="#">libtommath</a>	Public domain back-end with no licence restrictions.	Slower than <a href="#">GMP</a> .

### cpp\_int

```
#include <boost/multiprecision/cpp_int.hpp>
```

```

namespace boost{ namespace multiprecision{

typedef unspecified-type limb_type;

enum cpp_integer_type    { signed_magnitude, unsigned_magnitude };
enum cpp_int_check_type  { checked, unchecked };

template <unsigned MinDigits = 0,
          unsigned MaxDits = 0,
          cpp_integer_type SignType = signed_magnitude,
          cpp_int_check_type Checked = unchecked,
          class Allocator = std::allocator<limb_type> >
class cpp_int_backend;
//
// Expression templates default to et_off if there is no allocator:
//
template <unsigned MinDigits, unsigned MaxDigits, cpp_integer_type SignType,
cpp_int_check_type Checked>
struct expression_template_default<cpp_int_backend<MinDigits, MaxDigits, SignType, Checked, void> >
{ static const expression_template_option value = et_off; };

typedef number<cpp_int_backend<> >                cpp_int;    // arbitrary precision integer
typedef rational_adaptor<cpp_int_backend<> >      cpp_rational_backend;
typedef number<cpp_rational_backend>              cpp_rational; // arbitrary precision rational
number

// Fixed precision unsigned types:
typedef number<cpp_int_backend<128, 128, unsigned_magnitude, unchecked, void> >      uint128_t;
typedef number<cpp_int_backend<256, 256, unsigned_magnitude, unchecked, void> >      uint256_t;
typedef number<cpp_int_backend<512, 512, unsigned_magnitude, unchecked, void> >      uint512_t;
typedef number<cpp_int_backend<1024, 1024, unsigned_magnitude, unchecked, void> >     uint1024_t;

// Fixed precision signed types:
typedef number<cpp_int_backend<128, 128, signed_magnitude, unchecked, void> >         int128_t;
typedef number<cpp_int_backend<256, 256, signed_magnitude, unchecked, void> >         int256_t;
typedef number<cpp_int_backend<512, 512, signed_magnitude, unchecked, void> >         int512_t;
typedef number<cpp_int_backend<1024, 1024, signed_magnitude, unchecked, void> >        int1024_t;

// Over again, but with checking enabled this time:
typedef number<cpp_int_backend<0, 0, signed_magnitude, checked> >                    checked_cpp_int;
typedef rational_adaptor<cpp_int_backend<0, 0, signed_magnitude, checked> >         checked_cpp_rational_backend;
typedef number<cpp_rational_backend>                                                  checked_cpp_rational;

// Checked fixed precision unsigned types:
typedef number<cpp_int_backend<128, 128, unsigned_magnitude, checked, void> >         checked_uint128_t;
typedef number<cpp_int_backend<256, 256, unsigned_magnitude, checked, void> >         checked_uint256_t;
typedef number<cpp_int_backend<512, 512, unsigned_magnitude, checked, void> >         checked_uint512_t;
typedef number<cpp_int_backend<1024, 1024, unsigned_magnitude, checked, void> >        checked_uint1024_t;

// Fixed precision signed types:
typedef number<cpp_int_backend<128, 128, signed_magnitude, checked, void> >           checked_int128_t;
typedef number<cpp_int_backend<256, 256, signed_magnitude, checked, void> >           checked_int256_t;
typedef number<cpp_int_backend<512, 512, signed_magnitude, checked, void> >           checked_int512_t;
typedef number<cpp_int_backend<1024, 1024, signed_magnitude, checked, void> >          checked_int1024_t;

}} // namespaces

```

The `cpp_int_backend` type is normally used via one of the convenience typedefs given above.

This back-end is the "Swiss Army Knife" of integer types as it can represent both fixed and [arbitrary precision](#) integer types, and both signed and unsigned types. There are five template arguments:

MinBits	Determines the number of Bits to store directly within the object before resorting to dynamic memory allocation. When zero, this field is determined automatically based on how many bits can be stored in union with the dynamic storage header: setting a larger value may improve performance as larger integer values will be stored internally before memory allocation is required.
MaxBits	Determines the maximum number of bits to be stored in the type: resulting in a fixed precision type. When this value is the same as MinBits, then the Allocator parameter is ignored, as no dynamic memory allocation will ever be performed: in this situation the Allocator parameter should be set to type <code>void</code> . Note that this parameter should not be used simply to prevent large memory allocations, not only is that role better performed by the allocator, but fixed precision integers have a tendency to allocate all of MaxBits of storage more often than one would expect.
SignType	Determines whether the resulting type is signed or not. Note that for <a href="#">arbitrary precision</a> types this parameter must be <code>signed_magnitude</code> . For fixed precision types then this type may be either <code>signed_magnitude</code> or <code>unsigned_magnitude</code> .
Checked	This parameter has two values: <code>checked</code> or <code>unchecked</code> . See below.
Allocator	The allocator to use for dynamic memory allocation, or type <code>void</code> if <code>MaxBits == MinBits</code> .

When the template parameter `Checked` is set to `checked` then the result is a *checked-integer*, checked and unchecked integers have the following properties:

Condition	Checked-Integer	Unchecked-Integer
Numeric overflow in fixed precision arithmetic	Throws a <code>std::overflow_error</code> .	Performs arithmetic modulo $2^{\text{MaxBits}}$
Constructing an integer from a value that can not be represented in the target type	Throws a <code>std::range_error</code> .	Converts the value modulo $2^{\text{MaxBits}}$ , signed to unsigned conversions extract the last MaxBits bits of the 2's complement representation of the input value.
Unsigned subtraction yielding a negative value.	Throws a <code>std::range_error</code> .	Yields the value that would result from treating the unsigned type as a 2's complement signed type.
Attempting a bitwise operation on a negative value.	Throws a <code>std::range_error</code>	Yields the value, but not the bit pattern, that would result from performing the operation on a 2's complement integer type.

Things you should know when using this type:

- Default constructed `cpp_int_backends` have the value zero.
- Division by zero results in a `std::overflow_error` being thrown.
- Construction from a string that contains invalid non-numeric characters results in a `std::runtime_error` being thrown.
- Since the precision of `cpp_int_backend` is necessarily limited when the allocator parameter is `void`, care should be taken to avoid numeric overflow when using this type unless you actually want modulo-arithmetic behavior.
- The type uses a sign-magnitude representation internally, so type `int128_t` has 128-bits of precision plus an extra sign bit. In this respect the behaviour of these types differs from built-in 2's complement types. It might be tempting to use a 127-bit type

instead, and indeed this does work, but behaviour is still slightly different from a 2's complement built-in type as the min and max values are identical (apart from the sign), where as they differ by one for a true 2's complement type. That said it should be noted that there's no requirement for built-in types to be 2's complement either - it's simply that this is the most common format by far.

- Attempting to print negative values as either an Octal or Hexadecimal string results in a `std::runtime_error` being thrown, this is a direct consequence of the sign-magnitude representation.
- The fixed precision types `[checked_] [u]intXXX_t` have expression template support turned off - it seems to make little difference to the performance of these types either way - so we may as well have the faster compile times by turning the feature off.
- Unsigned types support subtraction - the result is "as if" a 2's complement operation had been performed as long as they are not *checked-integers* (see above). In other words they behave pretty much as a built in integer type would in this situation. So for example if we were using `uint128_t` then `uint128_t(1)-4` would result in the value `0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFD` of type `uint128_t`. However, had this operation been performed on `checked_uint128_t` then a `std::range_error` would have been thrown.
- Unary negation of unsigned types results in a compiler error (static assertion).
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.
- When used at fixed precision, the size of this type is always one machine word larger than you would expect for an N-bit integer: the extra word stores both the sign, and how many machine words in the integer are actually in use. The latter is an optimisation for larger fixed precision integers, so that a 1024-bit integer has almost the same performance characteristics as a 128-bit integer, rather than being 4 times slower for addition and 16 times slower for multiplication (assuming the values involved would always fit in 128 bits). Typically this means you can use an integer type wide enough for the "worst case scenario" with only minor performance degradation even if most of the time the arithmetic could in fact be done with a narrower type.
- When used at fixed precision and `MaxBits` is smaller than the number of bits in the largest native integer type, then internally `cpp_int_backend` switches to a "trivial" implementation where it is just a thin wrapper around a single integer. Note that it will still be slightly slower than a bare native integer, as it emulates a signed-magnitude representation rather than simply using the platforms native sign representation: this ensures there is no step change in behavior as a `cpp_int` grows in size.
- Fixed precision `cpp_int`'s have some support for `constexpr` values and user-defined literals, see [here](#) for the full description. For example `0xfffff_cppi1024` specifies a 1024-bit integer with the value `0xffff`. This can be used to generate compile time constants that are too large to fit into any built in number type.

### Example:

```
#include <boost/multiprecision/cpp_int.hpp>

using namespace boost::multiprecision;

int128_t v = 1;

// Do some fixed precision arithmetic:
for(unsigned i = 1; i <= 20; ++i)
    v *= i;

std::cout << v << std::endl; // prints 20!

// Repeat at arbitrary precision:
cpp_int u = 1;
for(unsigned i = 1; i <= 100; ++i)
    u *= i;

std::cout << u << std::endl; // prints 100!
```

### gmp\_int

```
#include <boost/multiprecision/gmp.hpp>
```

```
namespace boost{ namespace multiprecision{  
  
    class gmp_int;  
  
    typedef number<gmp_int >          mpz_int;  
  
}} // namespaces
```

The `gmp_int` back-end is used via the typedef `boost::multiprecision::mpz_int`. It acts as a thin wrapper around the [GMP](#) `mpz_t` to provide an integer type that is a drop-in replacement for the native C++ integer types, but with unlimited precision.

As well as the usual conversions from arithmetic and string types, type `mpz_int` is copy constructible and assignable from:

- The [GMP](#) native types: `mpf_t`, `mpz_t`, `mpq_t`.
- Instances of `number<T>` that are wrappers around those types: `number<gmp_float<N> >`, `number<gmp_rational>`.

It's also possible to access the underlying `mpz_t` via the `data()` member function of `gmp_int`.

Things you should know when using this type:

- No changes are made to the GMP library's global settings - so you can safely mix this type with existing code that uses [GMP](#).
- Default constructed `gmp_ints` have the value zero (this is GMP's default behavior).
- Formatted IO for this type does not support octal or hexadecimal notation for negative values, as a result performing formatted output on this type when the argument is negative and either of the flags `std::ios_base::oct` or `std::ios_base::hex` are set, will result in a `std::runtime_error` will be thrown.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid integer.
- Division by zero results in a `std::overflow_error` being thrown.
- Although this type is a wrapper around [GMP](#) it will work equally well with [MPIR](#). Indeed use of [MPIR](#) is recommended on Win32.
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.

### Example:

```
#include <boost/multiprecision/gmp.hpp>  
  
using namespace boost::multiprecision;  
  
mpz_int v = 1;  
  
// Do some arithmetic:  
for(unsigned i = 1; i <= 1000; ++i)  
    v *= i;  
  
std::cout << v << std::endl; // prints 1000!  
  
// Access the underlying representation:  
mpz_t z;  
mpz_init(z);  
mpz_set(z, v.backend().data());
```

## tom\_int

```
#include <boost/multiprecision/tommath.hpp>
```

```
namespace boost{ namespace multiprecision{  
  
    class tommath_int;  
  
    typedef number<tommath_int >          tom_int;  
  
}} // namespaces
```

The `tommath_int` back-end is used via the typedef `boost::multiprecision::tom_int`. It acts as a thin wrapper around the [libtommath](#) `tom_int` to provide an integer type that is a drop-in replacement for the native C++ integer types, but with unlimited precision.

Things you should know when using this type:

- Default constructed objects have the value zero (this is [libtommath](#)'s default behavior).
- Although `tom_int` is mostly a drop in replacement for the builtin integer types, it should be noted that it is a rather strange beast as it's a signed type that is not a 2's complement type. As a result the bitwise operations `|` `&` `^` will throw a `std::runtime_error` exception if either of the arguments is negative. Similarly the complement operator `~` is deliberately not implemented for this type.
- Formatted IO for this type does not support octal or hexadecimal notation for negative values, as a result performing formatted output on this type when the argument is negative and either of the flags `std::ios_base::oct` or `std::ios_base::hex` are set, will result in a `std::runtime_error` will be thrown.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid integer.
- Division by zero results in a `std::overflow_error` being thrown.

#### Example:

```
#include <boost/multiprecision/tommath.hpp>  
  
boost::multiprecision::tom_int v = 1;  
  
// Do some arithmetic:  
for(unsigned i = 1; i <= 1000; ++i)  
    v *= i;  
  
std::cout << v << std::endl; // prints 1000!  
std::cout << std::hex << v << std::endl; // prints 1000! in hex format  
  
try{  
    std::cout << std::hex << -v << std::endl; // Oops! can't print a negative value in hex format!  
}  
catch(const std::runtime_error& e)  
{  
    std::cout << e.what() << std::endl;  
}  
  
try{  
    // v is not a 2's complement type, bitwise operations are only supported  
    // on positive values:  
    v = -v & 2;  
}  
catch(const std::runtime_error& e)  
{  
    std::cout << e.what() << std::endl;  
}
```

## Examples

### Factorials

In this simple example, we'll write a routine to print out all of the factorials which will fit into a 128-bit integer. At the end of the routine we do some fancy iostream formatting of the results:

```
#include <boost/multiprecision/cpp_int.hpp>
#include <iostream>
#include <iomanip>
#include <vector>

void print_factorials()
{
    using boost::multiprecision::cpp_int;
    //
    // Print all the factorials that will fit inside a 128-bit integer.
    //
    // Begin by building a big table of factorials, once we know just how
    // large the largest is, we'll be able to "pretty format" the results.
    //
    // Calculate the largest number that will fit inside 128 bits, we could
    // also have used numeric_limits<int128_t>::max() for this value:
    cpp_int limit = (cpp_int(1) << 128) - 1;
    //
    // Our table of values:
    std::vector<cpp_int> results;
    //
    // Initial values:
    unsigned i = 1;
    cpp_int factorial = 1;
    //
    // Cycle through the factorials till we reach the limit:
    while(factorial < limit)
    {
        results.push_back(factorial);
        ++i;
        factorial *= i;
    }
    //
    // Lets see how many digits the largest factorial was:
    unsigned digits = results.back().str().size();
    //
    // Now print them out, using right justification, while we're at it
    // we'll indicate the limit of each integer type, so begin by defining
    // the limits for 16, 32, 64 etc bit integers:
    cpp_int limits[] = {
        (cpp_int(1) << 16) - 1,
        (cpp_int(1) << 32) - 1,
        (cpp_int(1) << 64) - 1,
        (cpp_int(1) << 128) - 1,
    };
    std::string bit_counts[] = { "16", "32", "64", "128" };
    unsigned current_limit = 0;
    for(unsigned j = 0; j < results.size(); ++j)
    {
        if(limits[current_limit] < results[j])
        {
            std::string message = "Limit of " + bit_counts[current_limit] + " bit integers";
            std::cout << std::setfill('.') << std::setw(digits+1) << std::right << message << std::setfill(' ');
        }
    }
}
```

```
fill(' ') << std::endl;
    ++current_limit;
}
std::cout << std::setw(digits + 1) << std::right << results[j] << std::endl;
}
```

The output from this routine is:

```

1
2
6
24
120
720
5040
40320
.....Limit of 16 bit integers
362880
3628800
39916800
479001600
.....Limit of 32 bit integers
6227020800
87178291200
1307674368000
20922789888000
355687428096000
6402373705728000
121645100408832000
2432902008176640000
.....Limit of 64 bit integers
51090942171709440000
1124000727777607680000
25852016738884976640000
620448401733239439360000
15511210043330985984000000
403291461126605635584000000
10888869450418352160768000000
304888344611713860501504000000
8841761993739701954543616000000
26525285981219105863630848000000
822283865417792281772556288000000
26313083693369353016721801216000000
868331761881188649551819440128000000
29523279903960414084761860964352000000
```

## Bit Operations

In this example we'll show how individual bits within an integer may be manipulated, we'll start with an often needed calculation of  $2^n - 1$ , which we could obviously implement like this:

```
using boost::multiprecision::cpp_int;

cpp_int b1(unsigned n)
{
    cpp_int r(1);
    return (r << n) - 1;
}
```



Calling:

```
std::cout << std::hex << std::showbase << b1(200) << std::endl;
```

Yields as expected:

```
0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```

However, we could equally just set the n'th bit in the result, like this:

```
cpp_int b2(unsigned n)
{
    cpp_int r(0);
    return --bit_set(r, n);
}
```

Note how the `bit_set` function sets the specified bit in its argument and then returns a reference to the result - which we can then simply decrement. The result from a call to `b2` is the same as that to `b1`.

We can equally test bits, so for example the n'th bit of the result returned from `b2` shouldn't be set unless we increment it first:

```
assert(!bit_test(b1(200), 200)); // OK
assert(bit_test(++b1(200), 200)); // OK
```

And of course if we flip the n'th bit after increment, then we should get back to zero:

```
assert(!bit_flip(++b1(200), 200)); // OK
```

## Floating Point Numbers

The following back-ends provide floating point arithmetic:

Backend Type	Header	Radix	Dependencies	Pros	Cons
<code>cpp_bin_float&lt;N&gt;</code>	<code>boost/multiprecision/cpp_bin_float.hpp</code>	2	None	Header only, all C++ implementation. Boost licence.	Approximately 2x slower than the <a href="#">MPFR</a> or <a href="#">GMP</a> libraries.
<code>cpp_dec_float&lt;N&gt;</code>	<code>boost/multiprecision/cpp_dec_float.hpp</code>	10	None	Header only, all C++ implementation. Boost licence.	Approximately 2x slower than the <a href="#">MPFR</a> or <a href="#">GMP</a> libraries.
<code>mpf_float&lt;N&gt;</code>	<code>boost/multiprecision/gmp.hpp</code>	2	<a href="#">GMP</a>	Very fast and efficient back-end.	Dependency on GNU licensed <a href="#">GMP</a> library.
<code>mpfr_float&lt;N&gt;</code>	<code>boost/multiprecision/mpfr.hpp</code>	2	<a href="#">GMP</a> and <a href="#">MPFR</a>	Very fast and efficient back-end, with its own standard library implementation.	Dependency on GNU licensed <a href="#">GMP</a> and <a href="#">MPFR</a> libraries.
<code>float128</code>	<code>boost/multiprecision/float128.hpp</code>	2	Either <a href="#">libquadmath</a> or the Intel C++ Math library.	Very fast and efficient back-end for 128-bit floating point values (113-bit mantissa, equivalent to FORTRAN's QUAD real)	Depends on the compiler being either recent GCC or Intel C++ versions.

## cpp\_bin\_float

```
#include <boost/multiprecision/cpp_bin_float.hpp>
```

```
namespace boost{ namespace multiprecision{

enum digit_base_type
{
    digit_base_2 = 2,
    digit_base_10 = 10
};

template <unsigned Digits, digit_base_type base = digit_base_10, class Allocator = void, class Exponent = int, ExponentMin = 0, ExponentMax = 0>
class cpp_bin_float;

typedef number<cpp_bin_float<50> > cpp_bin_float_50;
typedef number<cpp_bin_float<100> > cpp_bin_float_100;

typedef number<backends::cpp_bin_float<24, backends::digit_base_2, void, boost::int16_t, -126, 127>, et_off> cpp_bin_float_single;
typedef number<backends::cpp_bin_float<53, backends::digit_base_2, void, boost::int16_t, -1022, 1023>, et_off> cpp_bin_float_double;
typedef number<backends::cpp_bin_float<64, backends::digit_base_2, void, boost::int16_t, -16382, 16383>, et_off> cpp_bin_float_double_extended;
typedef number<backends::cpp_bin_float<113, backends::digit_base_2, void, boost::int16_t, -16382, 16383>, et_off> cpp_bin_float_quad;

}} // namespaces
```

The `cpp_bin_float` back-end is used in conjunction with `number`: It acts as an entirely C++ (header only and dependency free) floating-point number type that is a drop-in replacement for the native C++ floating-point types, but with much greater precision.

Type `cpp_bin_float` can be used at fixed precision by specifying a non-zero `Digits` template parameter. The typedefs `cpp_bin_float_50` and `cpp_bin_float_100` provide arithmetic types at 50 and 100 decimal digits precision respectively.

Optionally, you can specify whether the precision is specified in decimal digits or binary bits - for example to declare a `cpp_bin_float` with exactly the same precision as `double` one would use `number<cpp_bin_float<53, digit_base_2> >`. The typedefs `cpp_bin_float_single`, `cpp_bin_float_double`, `cpp_bin_float_quad` and `cpp_bin_float_double_extended` provide software analogues of the IEEE single, double and quad float data types, plus the Intel-extended-double type respectively. Note that while these types are functionally equivalent to the native IEEE types, but they do not have the same size or bit-layout as true IEEE compatible types.

Normally `cpp_bin_float` allocates no memory: all of the space required for its digits are allocated directly within the class. As a result care should be taken not to use the class with too high a digit count as stack space requirements can grow out of control. If that represents a problem then providing an allocator as a template parameter causes `cpp_bin_float` to dynamically allocate the memory it needs: this significantly reduces the size of `cpp_bin_float` and increases the viable upper limit on the number of digits at the expense of performance. However, please bear in mind that arithmetic operations rapidly become *very* expensive as the digit count grows: the current implementation really isn't optimized or designed for large digit counts. Note that since the actual type of the objects allocated is completely opaque, the suggestion would be to use an allocator with `void` value\_type, for example: `number<cpp_bin_float<1000, digit_base_10, std::allocator<void> > >`.

The final template parameters determine the type and range of the exponent: parameter `Exponent` can be any signed integer type, but note that `MinExponent` and `MaxExponent` can not go right up to the limits of the `Exponent` type as there has to be a little extra headroom for internal calculations. You will get a compile time error if this is the case. In addition if `MinExponent` or `MaxExponent` are zero, then the library will choose suitable values that are as large as possible given the constraints of the type and need for extra headroom for internal calculations.

There is full standard library and `numeric_limits` support available for this type.

Things you should know when using this type:

- Default constructed `cpp_bin_floats` have a value of zero.
- The radix of this type is 2, even when the precision is specified as decimal digits.

- The type supports both infinities and NaN's. An infinity is generated whenever the result would overflow, and a NaN is generated for any mathematically undefined operation.
- There is a `std::numeric_limits` specialisation for this type.
- Any number instantiated on this type, is convertible to any other number instantiated on this type - for example you can convert from `number<cpp_bin_float<50> >` to `number<cpp_bin_float<SomeOtherValue> >`. Narrowing conversions round to nearest and are explicit.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- All arithmetic operations are correctly rounded to nearest. String conversions and the `sqrt` function are also correctly rounded, but transcendental functions (`sin`, `cos`, `pow`, `exp` etc) are not.

### cpp\_bin\_float example:

```
#include <boost/multiprecision/cpp_bin_float.hpp>

using namespace boost::multiprecision;

// Operations at fixed precision and full numeric_limits support:
cpp_bin_float_100 b = 2;
std::cout << std::numeric_limits<cpp_bin_float_100>::digits << std::endl;
std::cout << std::numeric_limits<cpp_bin_float_100>::digits10 << std::endl;
// We can use any C++ std lib function, lets print all the digits as well:
std::cout << std::setprecision(std::numeric_limits<cpp_bin_float_100>::max_digits10)
    << log(b) << std::endl; // print log(2)
// We can also use any function from Boost.Math:
std::cout << boost::math::tgamma(b) << std::endl;
// These even work when the argument is an expression template:
std::cout << boost::math::tgamma(b * b) << std::endl;
// And since we have an extended exponent range we can generate some really large
// numbers here (4.0238726007709377354370243e+2564):
std::cout << boost::math::tgamma(cpp_bin_float_100(1000)) << std::endl;
```

### cpp\_dec\_float

```
#include <boost/multiprecision/cpp_dec_float.hpp>
```

```
namespace boost{ namespace multiprecision{

template <unsigned Digits10, class ExponentType = boost::int32_t, class Allocator = void>
class cpp_dec_float;

typedef number<cpp_dec_float<50> > cpp_dec_float_50;
typedef number<cpp_dec_float<100> > cpp_dec_float_100;

}} // namespaces
```

The `cpp_dec_float` back-end is used in conjunction with `number`: It acts as an entirely C++ (header only and dependency free) floating-point number type that is a drop-in replacement for the native C++ floating-point types, but with much greater precision.

Type `cpp_dec_float` can be used at fixed precision by specifying a non-zero `Digits10` template parameter. The typedefs `cpp_dec_float_50` and `cpp_dec_float_100` provide arithmetic types at 50 and 100 decimal digits precision respectively. Optionally, you can specify an integer type to use for the exponent, this defaults to a 32-bit integer type which is more than large enough for the vast majority of use cases, but larger types such as `long long` can also be specified if you need a truly huge exponent range. In any case the `ExponentType` must be a built in signed integer type at least 2 bytes and 16-bits wide.

Normally `cpp_dec_float` allocates no memory: all of the space required for its digits are allocated directly within the class. As a result care should be taken not to use the class with too high a digit count as stack space requirements can grow out of control. If that represents a problem then providing an allocator as the final template parameter causes `cpp_dec_float` to dynamically allocate the memory it needs: this significantly reduces the size of `cpp_dec_float` and increases the viable upper limit on the number of digits at the expense of performance. However, please bear in mind that arithmetic operations rapidly become *very* expensive as the digit count grows: the current implementation really isn't optimized or designed for large digit counts.

There is full standard library and `numeric_limits` support available for this type.

Things you should know when using this type:

- Default constructed `cpp_dec_floats` have a value of zero.
- The radix of this type is 10. As a result it can behave subtly differently from base-2 types.
- The type has a number of internal guard digits over and above those specified in the template argument. Normally these should not be visible to the user.
- The type supports both infinities and NaN's. An infinity is generated whenever the result would overflow, and a NaN is generated for any mathematically undefined operation.
- There is a `std::numeric_limits` specialisation for this type.
- Any number instantiated on this type, is convertible to any other number instantiated on this type - for example you can convert from `number<cpp_dec_float<50> >` to `number<cpp_dec_float<SomeOtherValue> >`. Narrowing conversions are truncating and explicit.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- The actual precision of a `cpp_dec_float` is always slightly higher than the number of digits specified in the template parameter, actually how much higher is an implementation detail but is always at least 8 decimal digits.
- Operations involving `cpp_dec_float` are always truncating. However, note that since their are guard digits in effect, in practice this has no real impact on accuracy for most use cases.

### cpp\_dec\_float example:

```
#include <boost/multiprecision/cpp_dec_float.hpp>

using namespace boost::multiprecision;

// Operations at fixed precision and full numeric_limits support:
cpp_dec_float_100 b = 2;
std::cout << std::numeric_limits<cpp_dec_float_100>::digits << std::endl;
// Note that digits10 is the same as digits, since we're base 10! :
std::cout << std::numeric_limits<cpp_dec_float_100>::digits10 << std::endl;
// We can use any C++ std lib function, lets print all the digits as well:
std::cout << std::setprecision(std::numeric_limits<cpp_dec_float_100>::max_digits10)
    << log(b) << std::endl; // print log(2)
// We can also use any function from Boost.Math:
std::cout << boost::math::tgamma(b) << std::endl;
// These even work when the argument is an expression template:
std::cout << boost::math::tgamma(b * b) << std::endl;
// And since we have an extended exponent range we can generate some really large
// numbers here (4.0238726007709377354370243e+2564):
std::cout << boost::math::tgamma(cpp_dec_float_100(1000)) << std::endl;
```

### gmp\_float

```
#include <boost/multiprecision/gmp.hpp>
```

```
namespace boost{ namespace multiprecision{  
  
    template <unsigned Digits10>  
    class gmp_float;  
  
    typedef number<gmp_float<50> >      mpf_float_50;  
    typedef number<gmp_float<100> >     mpf_float_100;  
    typedef number<gmp_float<500> >     mpf_float_500;  
    typedef number<gmp_float<1000> >    mpf_float_1000;  
    typedef number<gmp_float<0> >       mpf_float;  
  
}} // namespaces
```

The `gmp_float` back-end is used in conjunction with `number` : it acts as a thin wrapper around the [GMP](#) `mpf_t` to provide an real-number type that is a drop-in replacement for the native C++ floating-point types, but with much greater precision.

Type `gmp_float` can be used at fixed precision by specifying a non-zero `Digits10` template parameter, or at variable precision by setting the template argument to zero. The typedefs `mpf_float_50`, `mpf_float_100`, `mpf_float_500`, `mpf_float_1000` provide arithmetic types at 50, 100, 500 and 1000 decimal digits precision respectively. The typedef `mpf_float` provides a variable precision type whose precision can be controlled via the `numbers` member functions.



### Note

This type only provides standard library and `numeric_limits` support when the precision is fixed at compile time.

As well as the usual conversions from arithmetic and string types, instances of `number<mpf_float<N> >` are copy constructible and assignable from:

- The [GMP](#) native types `mpf_t`, `mpz_t`, `mpq_t`.
- The number wrappers around those types: `number<mpf_float<M> >`, `number<gmp_int>`, `number<gmp_rational>`.

It's also possible to access the underlying `mpf_t` via the `data()` member function of `gmp_float`.

Things you should know when using this type:

- Default constructed `gmp_floats` have the value zero (this is the [GMP](#) library's default behavior).
- No changes are made to the [GMP](#) library's global settings, so this type can be safely mixed with existing [GMP](#) code.
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.
- It is not possible to round-trip objects of this type to and from a string and get back exactly the same value. This appears to be a limitation of [GMP](#).
- Since the underlying [GMP](#) types have no notion of infinities or NaN's, care should be taken to avoid numeric overflow or division by zero. That latter will result in a `std::overflow_error` being thrown, while generating excessively large exponents may result in instability of the underlying [GMP](#) library (in testing, converting a number with an excessively large or small exponent to a string caused [GMP](#) to segfault).
- This type can equally be used with [MPIR](#) as the underlying implementation - indeed that is the recommended option on Win32.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- Division by zero results in a `std::overflow_error` being thrown.

**GMP example:**

```
#include <boost/multiprecision/gmp.hpp>

using namespace boost::multiprecision;

// Operations at variable precision and limited standard library support:
mpf_float a = 2;
mpf_float::default_precision(1000);
std::cout << mpf_float::default_precision() << std::endl;
std::cout << sqrt(a) << std::endl; // print root-2

// Operations at fixed precision and full standard library support:
mpf_float_100 b = 2;
std::cout << std::numeric_limits<mpf_float_100>::digits << std::endl;
// We can use any C++ std lib function:
std::cout << log(b) << std::endl; // print log(2)
// We can also use any function from Boost.Math:
std::cout << boost::math::tgamma(b) << std::endl;
// These even work when the argument is an expression template:
std::cout << boost::math::tgamma(b * b) << std::endl;

// Access the underlying representation:
mpf_t f;
mpf_init(f);
mpf_set(f, a.backend().data());
```

**mpfr\_float**

```
#include <boost/multiprecision/mpfr.hpp>
```

```
namespace boost{ namespace multiprecision{

enum mpfr_allocation_type
{
    allocate_stack,
    allocate_dynamic
};

template <unsigned Digits10, mpfr_allocation_type AllocateType = allocate_dynamic>
class mpfr_float_backend;

typedef number<mpfr_float_backend<50> >      mpfr_float_50;
typedef number<mpfr_float_backend<100> >     mpfr_float_100;
typedef number<mpfr_float_backend<500> >     mpfr_float_500;
typedef number<mpfr_float_backend<1000> >    mpfr_float_1000;
typedef number<mpfr_float_backend<0> >       mpfr_float;

typedef number<mpfr_float_backend<50, allocate_stack> >      static_mpfr_float_50;
typedef number<mpfr_float_backend<100, allocate_stack> >     static_mpfr_float_100;

}} // namespaces
```

The `mpfr_float_backend` type is used in conjunction with `number`: It acts as a thin wrapper around the [MPFR](#) `mpfr_t` to provide an real-number type that is a drop-in replacement for the native C++ floating-point types, but with much greater precision.

Type `mpfr_float_backend` can be used at fixed precision by specifying a non-zero `Digits10` template parameter, or at variable precision by setting the template argument to zero. The typedefs `mpfr_float_50`, `mpfr_float_100`, `mpfr_float_500`, `mpfr_float_1000` provide arithmetic types at 50, 100, 500 and 1000 decimal digits precision respectively. The typedef `mpfr_float` provides a variable precision type whose precision can be controlled via the `number`'s member functions.

In addition the second template parameter lets you choose between dynamic allocation (the default, and uses MPFR's normal allocation routines), or stack allocation (where all the memory required for the underlying data types is stored within `mpfr_float_backend`). The latter option can result in significantly faster code, at the expense of growing the size of `mpfr_float_backend`. It can only be used at fixed precision, and should only be used for lower digit counts. Note that we can not guarantee that using `allocate_stack` won't cause any calls to mpfr's allocation routines, as mpfr may call these inside it's own code. The following table gives an idea of the performance tradeoff's at 50 decimal digits precision<sup>2</sup>:

Type	Bessel function evaluation, relative times
<code>number&lt;mpfr_float_backend&lt;50, allocate_static&gt;, et_on&gt;</code>	1.0 (5.5s)
<code>number&lt;mpfr_float_backend&lt;50, allocate_static&gt;, et_off&gt;</code>	1.05 (5.8s)
<code>number&lt;mpfr_float_backend&lt;50, allocate_dynamic&gt;, et_on&gt;</code>	1.05 (5.8s)
<code>number&lt;mpfr_float_backend&lt;50, allocate_dynamic&gt;, et_off&gt;</code>	1.16 (6.4s)



### Note

This type only provides `numeric_limits` support when the precision is fixed at compile time.

As well as the usual conversions from arithmetic and string types, instances of `number<mpfr_float_backend<N> >` are copy constructible and assignable from:

- The [GMP](#) native types `mpf_t`, `mpz_t`, `mpq_t`.
- The [MPFR](#) native type `mpfr_t`.
- The `number` wrappers around those types: `number<mpfr_float_backend<M> >`, `number<mpf_float<M> >`, `number<gmp_int>`, `number<gmp_rational>`.

It's also possible to access the underlying `mpfr_t` via the `data()` member function of `mpfr_float_backend`.

Things you should know when using this type:

- A default constructed `mpfr_float_backend` is set to a NaN (this is the default [MPFR](#) behavior).
- All operations use round to nearest.
- No changes are made to [GMP](#) or [MPFR](#) global settings, so this type can coexist with existing [MPFR](#) or [GMP](#) code.
- The code can equally use [MPIR](#) in place of [GMP](#) - indeed that is the preferred option on Win32.
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- Division by zero results in an infinity.

<sup>2</sup> Compiled with VC++10 and /Ox, with MPFR-3.0.0 and MPIR-2.3.0



### MPFR example:

```
#include <boost/multiprecision/mpfr.hpp>

using namespace boost::multiprecision;

// Operations at variable precision and no numeric_limits support:
mpfr_float a = 2;
mpfr_float::default_precision(1000);
std::cout << mpfr_float::default_precision() << std::endl;
std::cout << sqrt(a) << std::endl; // print root-2

// Operations at fixed precision and full numeric_limits support:
mpfr_float_100 b = 2;
std::cout << std::numeric_limits<mpfr_float_100>::digits << std::endl;
// We can use any C++ std lib function:
std::cout << log(b) << std::endl; // print log(2)
// We can also use any function from Boost.Math:
std::cout << boost::math::tgamma(b) << std::endl;
// These even work when the argument is an expression template:
std::cout << boost::math::tgamma(b * b) << std::endl;

// Access the underlying data:
mpfr_t r;
mpfr_init(r);
mpfr_set(r, b.backend().data(), GMP_RNDN);
```

### float128

```
#include <boost/multiprecision/float128.hpp>
```

```
namespace boost{ namespace multiprecision{

class float128_backend;

typedef number<float128_backend, et_off>    float128;

}} // namespaces
```

The `float128` number type is a very thin wrapper around GCC's `float128` or Intel's `_Quad` data types and provides an real-number type that is a drop-in replacement for the native C++ floating-point types, but with a 113 bit mantissa, and compatible with FORTRAN's 128-bit QUAD real.

All the usual standard library and `numeric_limits` support are available, performance should be equivalent to the underlying native types: for example the LINPACK benchmarks for GCC's `float128` and `boost::multiprecision::float128` both achieved 5.6 MFLOPS<sup>3</sup>.

As well as the usual conversions from arithmetic and string types, instances of `float128` are copy constructible and assignable from GCC's `float128` and Intel's `_Quad` data types.

It's also possible to access the underlying `float128` or `_Quad` type via the `data()` member function of `float128_backend`.

Things you should know when using this type:

- Default constructed `float128`s have the value zero.
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.

---

<sup>3</sup> On 64-bit Ubuntu 11.10, GCC-4.8.0, Intel Core 2 Duo T5800.

- It is not possible to round-trip objects of this type to and from a string and get back exactly the same value when compiled with Intel's C++ compiler and using `_Quad` as the underlying type: this is a current limitation of our code. Round tripping when using `float128` as the underlying type is possible (both for GCC and Intel).
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- Division by zero results in an infinity being produced.
- Type `float128` can be used as a literal type (constexpr support).
- When using the Intel compiler, the underlying type defaults to `float128` if it's available and `_Quad` if not. You can override the default by defining either `BOOST_MP_USE_FLOAT128` or `BOOST_MP_USE_QUAD`.
- When the underlying type is Intel's `_Quad` type, the code must be compiled with the compiler option `-Qoption,cpp,--extended_float_type`.

### float128 example:

```
#include <boost/multiprecision/float128.hpp>

using namespace boost::multiprecision;

// Operations at 128-bit precision and full numeric_limits support:
float128 b = 2;
// There are 113-bits of precision:
std::cout << std::numeric_limits<float128>::digits << std::endl;
// Or 34 decimal places:
std::cout << std::numeric_limits<float128>::digits10 << std::endl;
// We can use any C++ std lib function, lets print all the digits as well:
std::cout << std::setprecision(std::numeric_limits<float128>::max_digits10)
    << log(b) << std::endl; // print log(2) = 0.693147180559945309417232121458176575
// We can also use any function from Boost.Math:
std::cout << boost::math::tgamma(b) << std::endl;
// And since we have an extended exponent range we can generate some really large
// numbers here (4.02387260077093773543702433923004111e+2564):
std::cout << boost::math::tgamma(float128(1000)) << std::endl;
//
// We can declare constants using GCC or Intel's native types, and the Q suffix,
// these can be declared constexpr if required:

constexpr float128 pi = 3.1415926535897932384626433832795028841971693993751058Q;
```

## Examples

### Area of Circle

Generic numeric programming employs templates to use the same code for different floating-point types and functions. Consider the area of a circle of radius  $r$ , given by

$$a = \pi * r^2$$

The area of a circle can be computed in generic programming using Boost.Math for the constant  $\pi$  as shown below:

```
#include <boost/math/constants/constants.hpp>

template<typename T>
inline T area_of_a_circle(T r)
{
    using boost::math::constants::pi;
    return pi<T>() * r * r;
}
```

It is possible to use `area_of_a_circle()` with built-in floating-point types as well as floating-point types from Boost.Multiprecision. In particular, consider a system with 4-byte single-precision float, 8-byte double-precision double and also the `cpp_dec_float_50` data type from Boost.Multiprecision with 50 decimal digits of precision.

We can compute and print the approximate area of a circle with radius 123/100 for `float`, `double` and `cpp_dec_float_50` with the program below.

```
#include <iostream>
#include <iomanip>
#include <boost/multiprecision/cpp_dec_float.hpp>

using boost::multiprecision::cpp_dec_float_50;

int main(int, char**)
{
    const float r_f(float(123) / 100);
    const float a_f = area_of_a_circle(r_f);

    const double r_d(double(123) / 100);
    const double a_d = area_of_a_circle(r_d);

    const cpp_dec_float_50 r_mp(cpp_dec_float_50(123) / 100);
    const cpp_dec_float_50 a_mp = area_of_a_circle(r_mp);

    // 4.75292
    std::cout
        << std::setprecision(std::numeric_limits<float>::digits10)
        << a_f
        << std::endl;

    // 4.752915525616
    std::cout
        << std::setprecision(std::numeric_limits<double>::digits10)
        << a_d
        << std::endl;

    // 4.7529155256159981904701331745635599135018975843146
    std::cout
        << std::setprecision(std::numeric_limits<cpp_dec_float_50>::digits10)
        << a_mp
        << std::endl;
}
```

In the next example we'll look at calling both standard library and Boost.Math functions from within generic code. We'll also show how to cope with template arguments which are expression-templates rather than number types.

## Defining a Special Function.

In this example we'll show several implementations of the [Jahnke and Emden Lambda function](#), each implementation a little more sophisticated than the last.

The Jahnke-Emden Lambda function is defined by the equation:

$$JahnkeEmden(v, z) = \Gamma(v+1) * J_v(z) / (z/2)^v$$

If we were to implement this at double precision using Boost.Math's facilities for the Gamma and Bessel function calls it would look like this:

```
double JEL1(double v, double z)
{
    return boost::math::tgamma(v + 1) * boost::math::cyl_bessel_j(v, z) / std::pow(z / 2, v);
}
```

Calling this function as:

```
std::cout << std::scientific << std::setprecision(std::numeric_limits<double>::digits10);
std::cout << JEL1(2.5, 0.5) << std::endl;
```

Yields the output:

```
9.822663964796047e-001
```

Now let's implement the function again, but this time using the multiprecision type `cpp_dec_float_50` as the argument type:

```
boost::multiprecision::cpp_dec_float_50
JEL2(boost::multiprecision::cpp_dec_float_50 v, boost::multiprecision::cpp_dec_float_50 z)
{
    return boost::math::tgamma(v + 1) * boost::math::cyl_bessel_j(v, z) / boost::multiprecision::pow(z / 2, v);
}
```

The implementation is almost the same as before, but with one key difference - we can no longer call `std::pow`, instead we must call the version inside the `boost::multiprecision` namespace. In point of fact, we could have omitted the namespace prefix on the call to `pow` since the right overload would have been found via [argument dependent lookup](#) in any case.

Note also that the first argument to `pow` along with the argument to `tgamma` in the above code are actually expression templates. The `pow` and `tgamma` functions will handle these arguments just fine.

Here's an example of how the function may be called:

```
std::cout << std::scientific << std::setprecision(std::numeric_limits<cpp_dec_float_50>::digits10);
std::cout << JEL2(cpp_dec_float_50(2.5), cpp_dec_float_50(0.5)) << std::endl;
```

Which outputs:

```
9.82266396479604757017335009796882833995903762577173e-01
```

Now that we've seen some non-template examples, let's repeat the code again, but this time as a template that can be called either with a builtin type (`float`, `double` etc), or with a multiprecision type:

```
template <class Float>
Float JEL3(Float v, Float z)
{
    using std::pow;
    return boost::math::tgamma(v + 1) * boost::math::cyl_bessel_j(v, z) / pow(z / 2, v);
}
```

Once again the code is almost the same as before, but the call to `pow` has changed yet again. We need the call to resolve to either `std::pow` (when the argument is a builtin type), or to `boost::multiprecision::pow` (when the argument is a multiprecision

type). We do that by making the call unqualified so that versions of `pow` defined in the same namespace as type `Float` are found via argument dependent lookup, while the `using std::pow` directive makes the standard library versions visible for builtin floating point types.

Let's call the function with both `double` and `multiprecision` arguments:

```
std::cout << std::scientific << std::setprecision(std::numeric_limits<double>::digits10);
std::cout << JEL3(2.5, 0.5) << std::endl;
std::cout << std::scientific << std::setprecision(std::numeric_limits<cpp_dec_float_50>::digits10);
std::cout << JEL3(cpp_dec_float_50(2.5), cpp_dec_float_50(0.5)) << std::endl;
```

Which outputs:

```
9.822663964796047e-001
9.82266396479604757017335009796882833995903762577173e-01
```

Unfortunately there is a problem with this version: if we were to call it like this:

```
boost::multiprecision::cpp_dec_float_50 v(2), z(0.5);
JEL3(v + 0.5, z);
```

Then we would get a long and inscrutable error message from the compiler: the problem here is that the first argument to `JEL3` is not a number type, but an expression template. We could obviously add a typecast to fix the issue:

```
JEL(cpp_dec_float_50(v + 0.5), z);
```

However, if we want the function `JEL` to be truly reusable, then a better solution might be preferred. To achieve this we can borrow some code from `Boost.Math` which calculates the return type of mixed-argument functions, here's how the new code looks now:

```
template <class Float1, class Float2>
typename boost::math::tools::promote_args<Float1, Float2>::type
    JEL4(Float1 v, Float2 z)
{
    using std::pow;
    return boost::math::tgamma(v + 1) * boost::math::cyl_bessel_j(v, z) / pow(z / 2, v);
}
```

As you can see the two arguments to the function are now separate template types, and the return type is computed using the `promote_args` metafunction from `Boost.Math`.

Now we can call:

```
std::cout << std::scientific << std::setprecision(std::numeric_limits<cpp_dec_float_100>::digits10);
std::cout << JEL4(cpp_dec_float_100(2) + 0.5, cpp_dec_float_100(0.5)) << std::endl;
```

And get 100 digits of output:

```
9.8226639647960475701733500979688283399590376257717309069410413822165082248153638454147004236848917775e-01
```

As a bonus, we can now call the function not just with expression templates, but with other mixed types as well: for example `float` and `double` or `int` and `double`, and the correct return type will be computed in each case.

Note that while in this case we didn't have to change the body of the function, in the general case any function like this which creates local variables internally would have to use `promote_args` to work out what type those variables should be, for example:

```
template <class Float1, class Float2>
typename boost::math::tools::promote_args<Float1, Float2>::type
JEL5(Float1 v, Float2 z)
{
    using std::pow;
    typedef typename boost::math::tools::promote_args<Float1, Float2>::type variable_type;
    variable_type t = pow(z / 2, v);
    return boost::math::tgamma(v + 1) * boost::math::cyl_bessel_j(v, z) / t;
}
```

## Calculating a Derivative

In this example we'll add even more power to generic numeric programming using not only different floating-point types but also function objects as template parameters. Consider some well-known central difference rules for numerically computing the first derivative of a function  $f'(x)$  with  $x \in \mathbb{R}$ :

$$\begin{aligned}
 f'(x) &\approx m_1 + O(dx^2) \\
 (1) \quad f'(x) &\approx \frac{4}{3}m_1 - \frac{1}{3}m_2 + O(dx^4) \\
 f'(x) &\approx \frac{3}{2}m_1 - \frac{3}{5}m_2 + \frac{1}{10}m_3 + O(dx^6)
 \end{aligned}$$

Where the difference terms  $m_n$  are given by:

$$\begin{aligned}
 m_1 &= \frac{f(x+dx) - f(x-dx)}{2dx} \\
 (2) \quad m_2 &= \frac{f(x+2dx) - f(x-2dx)}{4dx} \\
 m_3 &= \frac{f(x+3dx) - f(x-3dx)}{6dx}
 \end{aligned}$$

and  $dx$  is the step-size of the derivative.

The third formula in Equation 1 is a three-point central difference rule. It calculates the first derivative of  $f'(x)$  to  $O(dx^6)$ , where  $dx$  is the given step-size. For example, if the step-size is 0.01 this derivative calculation has about 6 decimal digits of precision - just about right for the 7 decimal digits of single-precision float. Let's make a generic template subroutine using this three-point central difference rule. In particular:

```
template<typename value_type, typename function_type>
value_type derivative(const value_type x, const value_type dx, function_type func)
{
    // Compute d/dx[func(*first)] using a three-point
    // central difference rule of O(dx^6).

    const value_type dx1 = dx;
    const value_type dx2 = dx1 * 2;
    const value_type dx3 = dx1 * 3;

    const value_type m1 = (func(x + dx1) - func(x - dx1)) / 2;
    const value_type m2 = (func(x + dx2) - func(x - dx2)) / 4;
    const value_type m3 = (func(x + dx3) - func(x - dx3)) / 6;

    const value_type fifteen_m1 = 15 * m1;
    const value_type six_m2      = 6 * m2;
    const value_type ten_dx1     = 10 * dx1;

    return ((fifteen_m1 - six_m2) + m3) / ten_dx1;
}
```

The `derivative()` template function can be used to compute the first derivative of any function to  $O(dx^6)$ . For example, consider the first derivative of  $\sin(x)$  evaluated at  $x = \pi/3$ . In other words,

$$(3) \quad \left. \frac{d}{dx} \sin x \right|_{x=\pi/3} = \cos \frac{\pi}{3} = \frac{1}{2}$$

The code below computes the derivative in Equation 3 for float, double and boost's multiple-precision type `cpp_dec_float_50`.

```
#include <iostream>
#include <iomanip>
#include <boost/multiprecision/cpp_dec_float.hpp>
#include <boost/math/constants/constants.hpp>

int main(int, char**)
{
    using boost::math::constants::pi;
    using boost::multiprecision::cpp_dec_float_50;
    //
    // We'll pass a function pointer for the function object passed to derivative,
    // the typecast is needed to select the correct overload of std::sin:
    //
    const float d_f = derivative(
        pi<float>() / 3,
        0.01F,
        static_cast<float(*)>(float)>(std::sin)
    );

    const double d_d = derivative(
        pi<double>() / 3,
        0.001,
        static_cast<double(*)>(double)>(std::sin)
    );
    //
    // In the cpp_dec_float_50 case, the sin function is multiply overloaded
    // to handle expression templates etc. As a result it's hard to take its
    // address without knowing about its implementation details. We'll use a
    // C++11 lambda expression to capture the call.
    // We also need a typecast on the first argument so we don't accidentally pass
    // an expression template to a template function:
    //
    const cpp_dec_float_50 d_mp = derivative(
        cpp_dec_float_50(pi<cpp_dec_float_50>() / 3),
        cpp_dec_float_50(1.0E-9),
        [](const cpp_dec_float_50& x) -> cpp_dec_float_50
        {
            return sin(x);
        }
    );

    // 5.000029e-001
    std::cout
        << std::setprecision(std::numeric_limits<float>::digits10)
        << d_f
        << std::endl;

    // 4.999999999998876e-001
    std::cout
        << std::setprecision(std::numeric_limits<double>::digits10)
        << d_d
        << std::endl;
```

We can take this a step further and use our derivative function to compute a partial derivative. For example if we take the incomplete gamma function  $P(a, z)$ , and take the derivative with respect to  $z$  at (2,2) then we can calculate the result as shown below, for good measure we'll compare with the "correct" result obtained from a call to `gamma_p_derivative`, the results agree to approximately 44 digits:

```
cpp_dec_float_50 gd = derivative(
    cpp_dec_float_50(2),
    cpp_dec_float_50(1.0E-9),
    [](const cpp_dec_float_50& x) ->cpp_dec_float_50
    {
        return boost::math::gamma_p(2, x);
    }
);
// 2.70670566473225383787998989944968806815263091819151e-01
std::cout
    << std::setprecision(std::numeric_limits<cpp_dec_float_50>::digits10)
    << gd
    << std::endl;
// 2.70670566473225383787998989944968806815253190143120e-01
std::cout << boost::math::gamma_p_derivative(
    cpp_dec_float_50(2), cpp_dec_float_50(2)) << std::endl;
```

Similar to the generic derivative example, we can calculate integrals in a similar manner:



```

template<typename value_type, typename function_type>
inline value_type integral(const value_type a,
                          const value_type b,
                          const value_type tol,
                          function_type func)
{
    unsigned n = 1U;

    value_type h = (b - a);
    value_type I = (func(a) + func(b)) * (h / 2);

    for(unsigned k = 0U; k < 8U; k++)
    {
        h /= 2;

        value_type sum(0);
        for(unsigned j = 1U; j <= n; j++)
        {
            sum += func(a + (value_type((j * 2) - 1) * h));
        }

        const value_type I0 = I;
        I = (I / 2) + (h * sum);

        const value_type ratio      = I0 / I;
        const value_type delta      = ratio - 1;
        const value_type delta_abs = ((delta < 0) ? -delta : delta);

        if((k > 1U) && (delta_abs < tol))
        {
            break;
        }

        n *= 2U;
    }

    return I;
}

```

The following sample program shows how the function can be called, we begin by defining a function object, which when integrated should yield the Bessel J function:

```

template<typename value_type>
class cyl_bessel_j_integral_rep
{
public:
    cyl_bessel_j_integral_rep(const unsigned N,
                             const value_type& X) : n(N), x(X) { }

    value_type operator()(const value_type& t) const
    {
        // pi * Jn(x) = Int_0^pi [cos(x * sin(t) - n*t) dt]
        return cos(x * sin(t) - (n * t));
    }

private:
    const unsigned n;
    const value_type x;
};

```

```
/* The function can now be called as follows: */
int main(int, char**)
{
    using boost::math::constants::pi;
    typedef boost::multiprecision::cpp_dec_float_50 mp_type;

    const float j2_f =
        integral(0.0F,
            pi<float>(),
            0.01F,
            cyl_bessel_j_integral_rep<float>(2U, 1.23F)) / pi<float>();

    const double j2_d =
        integral(0.0,
            pi<double>(),
            0.0001,
            cyl_bessel_j_integral_rep<double>(2U, 1.23)) / pi<double>();

    const mp_type j2_mp =
        integral(mp_type(0),
            pi<mp_type>(),
            mp_type(1.0E-20),
            cyl_bessel_j_integral_rep<mp_type>(2U, mp_type(123) / 100)) / pi<mp_type>();

    // 0.166369
    std::cout
        << std::setprecision(std::numeric_limits<float>::digits10)
        << j2_f
        << std::endl;

    // 0.166369383786814
    std::cout
        << std::setprecision(std::numeric_limits<double>::digits10)
        << j2_d
        << std::endl;

    // 0.16636938378681407351267852431513159437103348245333
    std::cout
        << std::setprecision(std::numeric_limits<mp_type>::digits10)
        << j2_mp
        << std::endl;

    //
    // Print true value for comparison:
    // 0.166369383786814073512678524315131594371033482453329
    std::cout << boost::math::cyl_bessel_j(2, mp_type(123) / 100) << std::endl;
}
```

## Polynomial Evaluation

In this example we'll look at polynomial evaluation, this is not only an important use case, but it's one that number performs particularly well at because the expression templates *completely eliminate all temporaries* from a [Horner polynomial evaluation scheme](#).

The following code evaluates  $\sin(x)$  as a polynomial, accurate to at least 64 decimal places:

```

using boost::multiprecision::cpp_dec_float;
typedef boost::multiprecision::number<cpp_dec_float<64> > mp_type;

mp_type mysin(const mp_type& x)
{
    // Approximation of sin(x * pi/2) for -1 <= x <= 1, using an order 63 polynomial.
    static const std::array<mp_type, 32U> coefs =
    {{
        ↵
mp_type("+1.5707963267948966192313216916397514420985846996875529104874722961539082031431044993140174126711"), //"),
        ↵
mp_type("-0.64596409750624625365575656389794573337969351178927307696134454382929989411386887578263960484"), // ↵
^3
        ↵
mp_type("+0.07969262624616704512050554949047802252091164235106119545663865720995702920146198554317279"), // ↵
^5
        ↵
mp_type("-0.0046817541353186881006854639339534378594950280185010575749538605102665157913157426229824"), // ↵
^7
        ↵
mp_type("+0.00016044118478735982187266087016347332970280754062061156858775174056686380286868007443"), // ↵
^9
        mp_type("-3.598843235212085340458540018208389404888495232432127661083907575106196374913134E-
6"), // ^11
        mp_type("+5.692172921967926811775255303592184372902829756054598109818158853197797542565E-
8"), // ^13
        mp_type("-6.688035109811467232478226335783138689956270985704278659373558497256423498E-
10"), // ^15
        mp_type("+6.066935731106195667101445665327140070166203261129845646380005577490472E-
12"), // ^17
        mp_type("-4.3770654673137422771842713137763190948628970300842263615764520034332E-14"), // ^19
        mp_type("+2.571422892860473866153865950420487369167895373255729246889168337E-16"), // ^21
        mp_type("-1.253899540535457665340073300390626396596970180355253776711660E-18"), // ^23
        mp_type("+5.1564551765802823395375998562329055050964428219501277474E-21"), // ^25
        mp_type("-1.812399312848887477410034071087545686586497030654642705E-23"), // ^27
        mp_type("+5.50728578652238583570585513920522536675023562254864E-26"), // ^29
        mp_type("-1.461148710664467988723468673933026649943084902958E-28"), // ^31
        mp_type("+3.41405297003316172502972039913417222912445427E-31"), // ^33
        mp_type("-7.07885550810745570069916712806856538290251E-34"), // ^35
        mp_type("+1.31128947968267628970845439024155655665E-36"), // ^37
        mp_type("-2.18318293181145698535113946654065918E-39"), // ^39
        mp_type("+3.28462680978498856345937578502923E-42"), // ^41
        mp_type("-4.48753699028101089490067137298E-45"), // ^43
        mp_type("+5.59219884208696457859353716E-48"), // ^45
        mp_type("-6.38214503973500471720565E-51"), // ^47
        mp_type("+6.69528558381794452556E-54"), // ^49
        mp_type("-6.47841373182350206E-57"), // ^51
        mp_type("+5.800016389666445E-60"), // ^53
        mp_type("-4.818507347289E-63"), // ^55
        mp_type("+3.724683686E-66"), // ^57
        mp_type("-2.6856479E-69"), // ^59
        mp_type("+1.81046E-72"), // ^61
        mp_type("-1.133E-75"), // ^63
    }};

    const mp_type v = x * 2 / boost::math::constants::pi<mp_type>();
    const mp_type x2 = (v * v);
    //
    // Polynomial evaluation follows, if mp_type allocates memory then
    // just one such allocation occurs - to initialize the variable "sum" -
    // and no temporaries are created at all.
    //
    const mp_type sum = ((((((((((((((((((((((((((((((((((((((( + coefs[31U]

```

```
* x2 + coefs[30U])
* x2 + coefs[29U])
* x2 + coefs[28U])
* x2 + coefs[27U])
* x2 + coefs[26U])
* x2 + coefs[25U])
* x2 + coefs[24U])
* x2 + coefs[23U])
* x2 + coefs[22U])
* x2 + coefs[21U])
* x2 + coefs[20U])
* x2 + coefs[19U])
* x2 + coefs[18U])
* x2 + coefs[17U])
* x2 + coefs[16U])
* x2 + coefs[15U])
* x2 + coefs[14U])
* x2 + coefs[13U])
* x2 + coefs[12U])
* x2 + coefs[11U])
* x2 + coefs[10U])
* x2 + coefs[9U])
* x2 + coefs[8U])
* x2 + coefs[7U])
* x2 + coefs[6U])
* x2 + coefs[5U])
* x2 + coefs[4U])
* x2 + coefs[3U])
* x2 + coefs[2U])
* x2 + coefs[1U])
* x2 + coefs[0U])
* v;

return sum;
}
```

Calling the function like so:

```
mp_type pid4 = boost::math::constants::pi<mp_type>() / 4;
std::cout << std::setprecision(std::numeric_limits< ::mp_type>::digits10) << std::scientific;
std::cout << mysin(pid4) << std::endl;
```

Yields the expected output:

```
7.0710678118654752440084436210484903928483593768847403658833986900e-01
```

## Interval Number Types

There is one currently only one interval number type supported - [MPFI](#).

### `mpfi_float`

```
#include <boost/multiprecision/mpfi.hpp>
```

```

namespace boost{ namespace multiprecision{

template <unsigned Digits10>
class mpfi_float_backend;

typedef number<mpfi_float_backend<50> >      mpfi_float_50;
typedef number<mpfi_float_backend<100> >     mpfi_float_100;
typedef number<mpfi_float_backend<500> >     mpfi_float_500;
typedef number<mpfi_float_backend<1000> >    mpfi_float_1000;
typedef number<mpfi_float_backend<0> >      mpfi_float;

}} // namespaces

```

The `mpfi_float_backend` type is used in conjunction with `number`: It acts as a thin wrapper around the [MPFI](#) `mpfi_t` to provide an real-number type that is a drop-in replacement for the native C++ floating-point types, but with much greater precision and implementing interval arithmetic.

Type `mpfi_float_backend` can be used at fixed precision by specifying a non-zero `Digits10` template parameter, or at variable precision by setting the template argument to zero. The typedefs `mpfi_float_50`, `mpfi_float_100`, `mpfi_float_500`, `mpfi_float_1000` provide arithmetic types at 50, 100, 500 and 1000 decimal digits precision respectively. The typedef `mpfi_float` provides a variable precision type whose precision can be controlled via the `numbers` member functions.



### Note

This type only provides `numeric_limits` support when the precision is fixed at compile time.

As well as the usual conversions from arithmetic and string types, instances of `number<mpfi_float_backend<N> >` are copy constructible and assignable from:

- The [MPFI](#) native type `mpfi_t`.
- The number wrappers around [MPFI](#) or [MPFR](#): `number<mpfi_float_backend<M> >` and `number<mpfr_float<M> >`.
- There is a two argument constructor taking two `number<mpfr_float<M> >` arguments specifying the interval.

It's also possible to access the underlying `mpfi_t` via the `data()` member function of `mpfi_float_backend`.

Things you should know when using this type:

- A default constructed `mpfi_float_backend` is set to a NaN (this is the default [MPFI](#) behavior).
- No changes are made to [GMP](#) or [MPFR](#) global settings, so this type can coexist with existing [MPFR](#) or [GMP](#) code.
- The code can equally use [MPIR](#) in place of [GMP](#) - indeed that is the preferred option on Win32.
- This backend supports rvalue-references and is move-aware, making instantiations of `number` on this backend move aware.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid floating point number.
- Division by zero results in an infinity.

There are some additional non member functions for working on intervals:

```

template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfr_float_backend<Digits10>, ExpressionTemplates> lower(const number<mpfi_float_backend<Di-
digits10>, ExpressionTemplates>& val);

```

Returns the lower end of the interval.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfr_float_backend<Digits10>, ExpressionTemplates> upper(const number<mpfi_float_backend<Di-
digits10>, ExpressionTemplates>& val);
```

Returns the upper end of the interval.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfr_float_backend<Digits10>, ExpressionTemplates> median(const number<mpfi_float_backend<Di-
digits10>, ExpressionTemplates>& val);
```

Returns the mid point of the interval.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfr_float_backend<Digits10>, ExpressionTemplates> width(const number<mpfi_float_backend<Di-
digits10>, ExpressionTemplates>& val);
```

Returns the absolute width of the interval.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfi_float_backend<Digits10>, ExpressionTemplates> intersect(
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a,
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& b);
```

Returns the interval which is the intersection of the  $a$  and  $b$ . Returns an unspecified empty interval if there is no such intersection.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
number<mpfi_float_backend<Digits10>, ExpressionTemplates> hull(
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a,
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& b);
```

Returns the interval which is the union of  $a$  and  $b$ .

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool overlap(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a,
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& b);
```

Returns true only if the intervals  $a$  and  $b$  overlap.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates1, expression_tem-
plate_option ExpressionTemplates2>
bool in(const number<mpfr_float_backend<Digits10>, ExpressionTemplates1>& a,
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates2>& b);
```

Returns true only if point  $a$  is contained within the interval  $b$ .

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool zero_in(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a);
```

Returns true only if the interval  $a$  contains the value zero.

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool subset(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a,
    const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& b);
```

Returns true only if  $a$  is a subset of  $b$ .

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool proper_subset(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a,
                  const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& b);
```

Returns true only if  $a$  is a proper subset of  $b$ .

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool empty(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a);
```

Returns true only if  $a$  is an empty interval, equivalent to  $\text{upper}(a) < \text{lower}(a)$ .

```
template <unsigned Digits10, expression_template_option ExpressionTemplates>
bool singleton(const number<mpfi_float_backend<Digits10>, ExpressionTemplates>& a);
```

Returns true if  $\text{lower}(a) == \text{upper}(a)$ .

### MPFI example:

```
#include <boost/multiprecision/mpfi.hpp>

using namespace boost::multiprecision;

// Operations at variable precision and no numeric_limits support:
mpfi_float a = 2;
mpfi_float::default_precision(1000);
std::cout << mpfi_float::default_precision() << std::endl;
std::cout << sqrt(a) << std::endl; // print root-2

// Operations at fixed precision and full numeric_limits support:
mpfi_float_100 b = 2;
std::cout << std::numeric_limits<mpfi_float_100>::digits << std::endl;
// We can use any C++ std lib function:
std::cout << log(b) << std::endl; // print log(2)

// Access the underlying data:
mpfi_t r;
mpfi_init(r);
mpfi_set(r, b.backend().data());

// Construct some explicit intervals and perform set operations:
mpfi_float_50 i1(1, 2), i2(1.5, 2.5);
std::cout << intersect(i1, i2) << std::endl;
std::cout << hull(i1, i2) << std::endl;
std::cout << overlap(i1, i2) << std::endl;
std::cout << subset(i1, i2) << std::endl;
```

## Rational Number Types

The following back-ends provide rational number arithmetic:

Backend Type	Header	Radix	Dependencies	Pros	Cons
cpp_rational	boost/multiprecision/cpp_int.hpp	2	None	An all C++ Boost-licensed implementation.	Slower than <a href="#">GMP</a> .
gmp_rational	boost/multiprecision/gmp.hpp	2	<a href="#">GMP</a>	Very fast and efficient back-end.	Dependency on GNU licensed <a href="#">GMP</a> library.
tommath_rational	boost/multiprecision/tommath.hpp	2	<a href="#">libtommath</a>	All C/C++ implementation that's Boost Software Licence compatible.	Slower than <a href="#">GMP</a> .
rational_adaptor	boost/multiprecision/rational_adaptor.hpp	N/A	none	All C++ adaptor that allows any integer back-end type to be used as a rational type.	Requires an underlying integer back-end type.
boost::rational	boost/rational.hpp	N/A	None	A C++ rational number type that can be used with any integer type.	The expression templates used by number end up being "hidden" inside boost::rational: performance may well suffer as a result.

## cpp\_rational

```
#include <boost/multiprecision/cpp_int.hpp>
```

```
namespace boost{ namespace multiprecision{
    typedef rational_adaptor<cpp_int_backend<> >      cpp_rational_backend;
    typedef number<cpp_rational_backend>              cpp_rational;
}} // namespaces
```

The `cpp_rational_backend` type is used via the typedef `boost::multiprecision::cpp_rational`. It provides a rational number type that is a drop-in replacement for the native C++ number types, but with unlimited precision.

As well as the usual conversions from arithmetic and string types, instances of `cpp_rational` are copy constructible and assignable from type `cpp_int`.

There is also a two argument constructor that accepts a numerator and denominator: both of type `cpp_int`.

There are also non-member functions:

```
cpp_int numerator(const cpp_rational&);
cpp_int denominator(const cpp_rational&);
```

which return the numerator and denominator of the number.



Things you should know when using this type:

- Default constructed `cpp_rationals` have the value zero.
- Division by zero results in a `std::overflow_error` being thrown.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid rational number.

#### Example:

```
#include <boost/multiprecision/cpp_int.hpp>

using namespace boost::multiprecision;

cpp_rational v = 1;

// Do some arithmetic:
for(unsigned i = 1; i <= 1000; ++i)
    v *= i;
v /= 10;

std::cout << v << std::endl; // prints 1000! / 10
std::cout << numerator(v) << std::endl;
std::cout << denominator(v) << std::endl;

cpp_rational w(2, 3); // component wise constructor
std::cout << w << std::endl; // prints 2/3
```

## gmp\_rational

```
#include <boost/multiprecision/gmp.hpp>
```

```
namespace boost{ namespace multiprecision{

class gmp_rational;

typedef number<gmp_rational >          mpq_rational;

}} // namespaces
```

The `gmp_rational` back-end is used via the typedef `boost::multiprecision::mpq_rational`. It acts as a thin wrapper around the [GMP](#) `mpq_t` to provide a rational number type that is a drop-in replacement for the native C++ number types, but with unlimited precision.

As well as the usual conversions from arithmetic and string types, instances of `number<gmp_rational>` are copy constructible and assignable from:

- The [GMP](#) native types: `mpz_t`, `mpq_t`.
- `number<gmp_int>`.

There is also a two-argument constructor that accepts a numerator and denominator (both of type `number<gmp_int>`).

There are also non-member functions:

```
mpz_int numerator(const mpq_rational&);
mpz_int denominator(const mpq_rational&);
```

which return the numerator and denominator of the number.

It's also possible to access the underlying `mpq_t` via the `data()` member function of `mpq_rational`.

Things you should know when using this type:

- Default constructed `mpq_rationals` have the value zero (this is the [GMP](#) default behavior).
- Division by zero results in a `std::overflow_error` being thrown.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid rational number.
- No changes are made to the [GMP](#) library's global settings, so this type can coexist with existing [GMP](#) code.
- The code can equally be used with [MPIR](#) as the underlying library - indeed that is the preferred option on Win32.

### Example:

```
#include <boost/multiprecision/gmp.hpp>

using namespace boost::multiprecision;

mpq_rational v = 1;

// Do some arithmetic:
for(unsigned i = 1; i <= 1000; ++i)
    v *= i;
v /= 10;

std::cout << v << std::endl; // prints 1000! / 10
std::cout << numerator(v) << std::endl;
std::cout << denominator(v) << std::endl;

mpq_rational w(2, 3); // component wise constructor
std::cout << w << std::endl; // prints 2/3

// Access the underlying data:
mpq_t q;
mpq_init(q);
mpq_set(q, v.backend().data());
```

## tommath\_rational

```
#include <boost/multiprecision/tommath.hpp>
```

```
namespace boost{ namespace multiprecision{

typedef rational_adpater<tommath_int>          tommath_rational;
typedef number<tommath_rational >             tom_rational;

}} // namespaces
```

The `tommath_rational` back-end is used via the typedef `boost::multiprecision::tom_rational`. It acts as a thin wrapper around `boost::rational<tom_int>` to provide a rational number type that is a drop-in replacement for the native C++ number types, but with unlimited precision.

The advantage of using this type rather than `boost::rational<tom_int>` directly, is that it is expression-template enabled, greatly reducing the number of temporaries created in complex expressions.

There are also non-member functions:

```
tom_int numerator(const tom_rational&);  
tom_int denominator(const tom_rational&);
```

which return the numerator and denominator of the number.

Things you should know when using this type:

- Default constructed `tom_rationals` have the value zero (this the inherited `Boost.Rational` behavior).
- Division by zero results in a `std::overflow_error` being thrown.
- Conversion from a string results in a `std::runtime_error` being thrown if the string can not be interpreted as a valid rational number.
- No changes are made to `libtommath`'s global state, so this type can safely coexist with other `libtommath` code.
- Performance of this type has been found to be pretty poor - this need further investigation - but it appears that `Boost.Rational` needs some improvement in this area.

### Example:

```
#include <boost/multiprecision/tommath.hpp>  
  
using namespace boost::multiprecision;  
  
tom_rational v = 1;  
  
// Do some arithmetic:  
for(unsigned i = 1; i <= 1000; ++i)  
    v *= i;  
v /= 10;  
  
std::cout << v << std::endl; // prints 1000! / 10  
std::cout << numerator(v) << std::endl;  
std::cout << denominator(v) << std::endl;  
  
tom_rational w(2, 3); // Component wise constructor  
std::cout << w << std::endl; // prints 2/3
```

## Use With Boost.Rational

All of the integer types in this library can be used as template arguments to `boost::rational<IntType>`.

Note that using the library in this way largely negates the effect of the expression templates in number.

## rational\_adaptor

```
namespace boost{ namespace multiprecision{  
  
    template <class IntBackend>  
    class rational_adpater;  
  
}}
```

The class template `rational_adaptor` is a back-end for number which converts any existing integer back-end into a rational-number back-end.

So for example, given an integer back-end type `MyIntegerBackend`, the use would be something like:

```
typedef number<MyIntegerBackend> MyInt;
typedef number<rational_adaptor<MyIntegerBackend> > MyRational;

MyRational r = 2;
r /= 3;
MyInt i = numerator(r);
assert(i == 2);
```

## Miscellaneous Number Types.

Backend types listed in this section are predominantly designed to aid debugging.

### logged\_adaptor

```
#include <boost/multiprecision/logged_adaptor.hpp>
```

```
namespace boost{ namespace multiprecision{

template <class Backend>
void log_postfix_event(const Backend& result, const char* event_description);
template <class Backend, class T>
void log_postfix_event(const Backend& result1, const T& result2, const char* event_description);

template <class Backend>
void log_prefix_event(const Backend& arg1, const char* event_description);
template <class Backend, class T>
void log_prefix_event(const Backend& arg1, const T& arg2, const char* event_description);
template <class Backend, class T, class U>
void log_prefix_event(const Backend& arg1, const T& arg2, const U& arg3, const char* event_description);
template <class Backend, class T, class U, class V>
void log_prefix_event(const Backend& arg1, const T& arg2, const U& arg3, const V& arg4, const char* event_description);

template <Backend>
class logged_adaptor;

}} // namespaces
```

The `logged_adaptor` type is used in conjunction with `number` and some other backend type: it acts as a thin wrapper around some other backend to class `number` and logs all the events that take place on that object. Before any number operation takes place, it calls `log_prefix_event` with the arguments to the operation (up to 4), plus a string describing the operation. Then after the operation it calls `log_postfix_event` with the result of the operation, plus a string describing the operation. Optionally, `log_postfix_event` takes a second result argument: this occurs when the result of the operation is not a number, for example when `fpclassify` is called, `log_postfix_event` will be called with `result1` being the argument to the function, and `result2` being the integer result of `fpclassify`.

The default versions of `log_prefix_event` and `log_postfix_event` do nothing, it is therefore up to the user to overload these for the particular backend being observed.

This type provides `numeric_limits` support whenever the template argument `Backend` does so.

This type is particularly useful when combined with an interval number type - in this case we can use `log_postfix_event` to monitor the error accumulated after each operation. We could either set some kind of trap whenever the accumulated error exceeds some threshold, or simply print out diagnostic information. Using this technique we can quickly locate the cause of numerical instability in a particular routine. The following example demonstrates this technique in a trivial algorithm that deliberately introduces cancellation error:

```
#include <boost/multiprecision/mpfi.hpp>
#include <boost/multiprecision/logged_adaptor.hpp>
#include <iostream>
#include <iomanip>
//
// Begin by overloading log_postfix_event so we can capture each arithmetic event as it happens:
//
namespace boost{ namespace multiprecision{

template <unsigned D>
inline void log_postfix_event(const mpfi_float_backend<D>& val, const char* event_description)
{
    // Print out the (relative) diameter of the interval:
    using namespace boost::multiprecision;
    number<mpfr_float_backend<D> > diam;
    mpfi_diam(diam.backend().data(), val.data());
    std::cout << "Diameter was " << diam << " after operation: " << event_description << std::endl;
}

template <unsigned D, class T>
inline void log_postfix_event(const mpfi_float_backend<D>&, const T&, const char* event_description)
{
    // This version is never called in this example.
}

}}

int main()
{
    using namespace boost::multiprecision;
    typedef number<logged_adaptor<mpfi_float_backend<17> > > logged_type;
    //
    // Test case deliberately introduces cancellation error, relative size of interval
    // gradually gets larger after each operation:
    //
    logged_type a = 1;
    a /= 10;

    for(unsigned i = 0; i < 13; ++i)
    {
        logged_type b = a * 9;
        b /= 10;
        a -= b;
    }
    std::cout << "Final value was: " << a << std::endl;
    return 0;
}
```

When we examine program output we can clearly see that the diameter of the interval increases after each subtraction:

```
Diameter was nan after operation: Default construct
Diameter was 0 after operation: Assignment from arithmetic type
Diameter was 4.33681e-18 after operation: /=
Diameter was nan after operation: Default construct
Diameter was 7.70988e-18 after operation: *
Diameter was 9.63735e-18 after operation: /=
Diameter was 1.30104e-16 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 1.30104e-16 after operation: *
Diameter was 1.38537e-16 after operation: /=
Diameter was 2.54788e-15 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 2.54788e-15 after operation: *
Diameter was 2.54863e-15 after operation: /=
Diameter was 4.84164e-14 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 4.84164e-14 after operation: *
Diameter was 4.84221e-14 after operation: /=
Diameter was 9.19962e-13 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 9.19962e-13 after operation: *
Diameter was 9.19966e-13 after operation: /=
Diameter was 1.74793e-11 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 1.74793e-11 after operation: *
Diameter was 1.74793e-11 after operation: /=
Diameter was 3.32107e-10 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 3.32107e-10 after operation: *
Diameter was 3.32107e-10 after operation: /=
Diameter was 6.31003e-09 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 6.31003e-09 after operation: *
Diameter was 6.31003e-09 after operation: /=
Diameter was 1.19891e-07 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 1.19891e-07 after operation: *
Diameter was 1.19891e-07 after operation: /=
Diameter was 2.27792e-06 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 2.27792e-06 after operation: *
Diameter was 2.27792e-06 after operation: /=
Diameter was 4.32805e-05 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 4.32805e-05 after operation: *
Diameter was 4.32805e-05 after operation: /=
Diameter was 0.00082233 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 0.00082233 after operation: *
Diameter was 0.00082233 after operation: /=
Diameter was 0.0156243 after operation: -=
Diameter was nan after operation: Default construct
Diameter was 0.0156243 after operation: *
Diameter was 0.0156243 after operation: /=
Diameter was 0.296861 after operation: -=
Final value was: {8.51569e-15,1.14843e-14}
```

## debug\_adaptor

```
#include <boost/multiprecision/debug_adaptor.hpp>
```

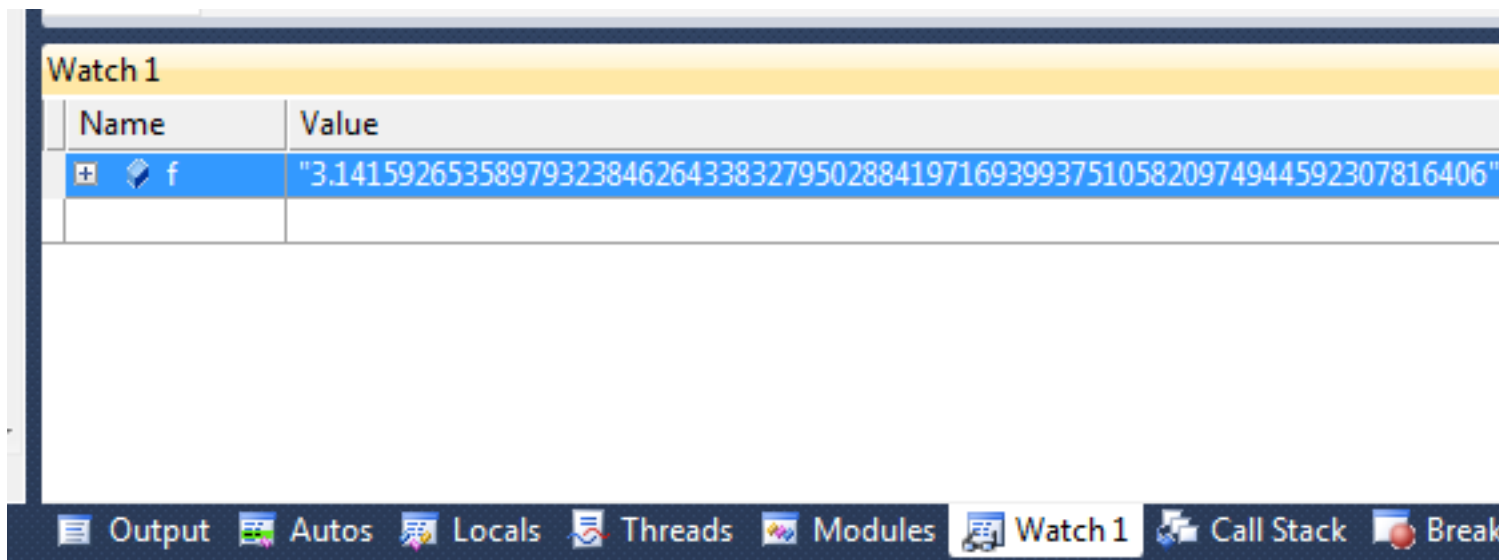
```
namespace boost{ namespace multiprecision{  
  
    template <Backend>  
    class debug_adaptor;  
  
}} // namespaces
```

The `debug_adaptor` type is used in conjunction with `number` and some other backend type: it acts as a thin wrapper around some other backend to class `number` and intercepts all operations on that object storing the result as a string within itself.

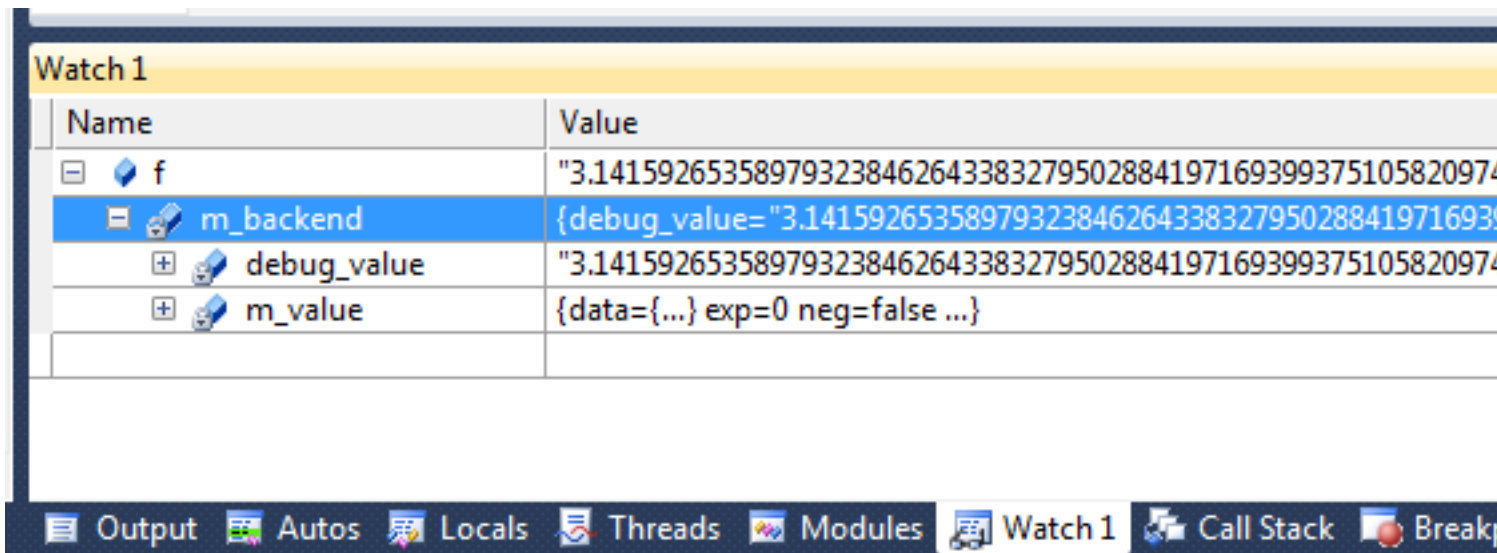
This type provides `numeric_limits` support whenever the template argument `Backend` does so.

This type is particularly useful when your debugger provides a good view of `std::string`: when this is the case multiprecision values can easily be inspected in the debugger by looking at the `debug_value` member of `debug_adaptor`. The down side of this approach is that runtimes are much slower when using this type. Set against that it can make debugging very much easier, certainly much easier than sprinkling code with `printf` statements.

When used in conjunction with the Visual C++ debugger visualisers, the value of a multiprecision type that uses this backend is displayed in the debugger just a builtin value would be, here we're inspecting a value of type `number<debug_adaptor<cpp_dec_float<50> > >`:

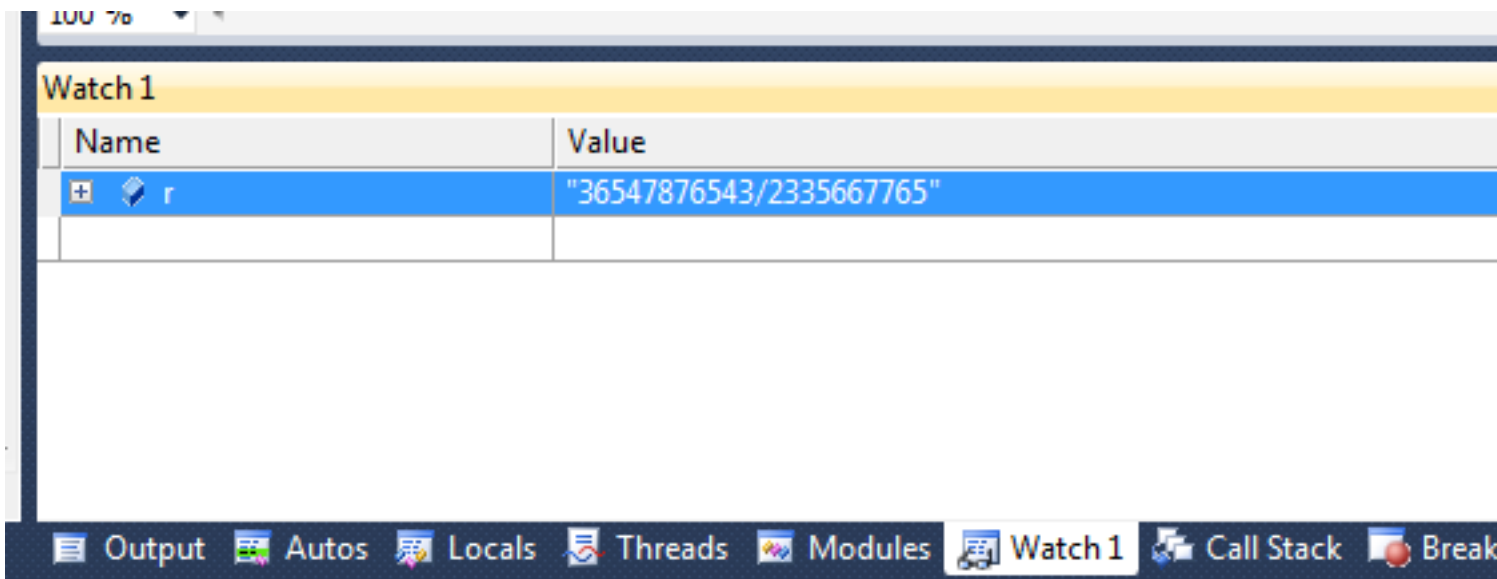


Otherwise you will need to expand out the view and look at the "debug\_value" member:



Name	Value
f	"3.141592653589793238462643383279502884197169399375105820974"
m_backend	{debug_value="3.141592653589793238462643383279502884197169399375105820974"
debug_value	"3.141592653589793238462643383279502884197169399375105820974"
m_value	{data={...} exp=0 neg=false ...}

It works for all the backend types equally too, here it is inspecting a number<debug\_adaptor<gmp\_rational> >:



Name	Value
r	"36547876543/2335667765"

## Visual C++ Debugger Visualizers

Let's face it debugger multiprecision numbers is hard - simply because we can't easily inspect the value of the numbers. Visual C++ provides a partial solution in the shape of "visualizers" which provide improved views of complex data structures, these visualizers need to be added to the [Visualizer] section of `autoexp.dat` located in the `Common7/Packages/Debugger` directory of your Visual Studio installation. The actual visualizer code is in the sandbox [here](#) - just cut and paste the code into your `autoexp.dat` file.

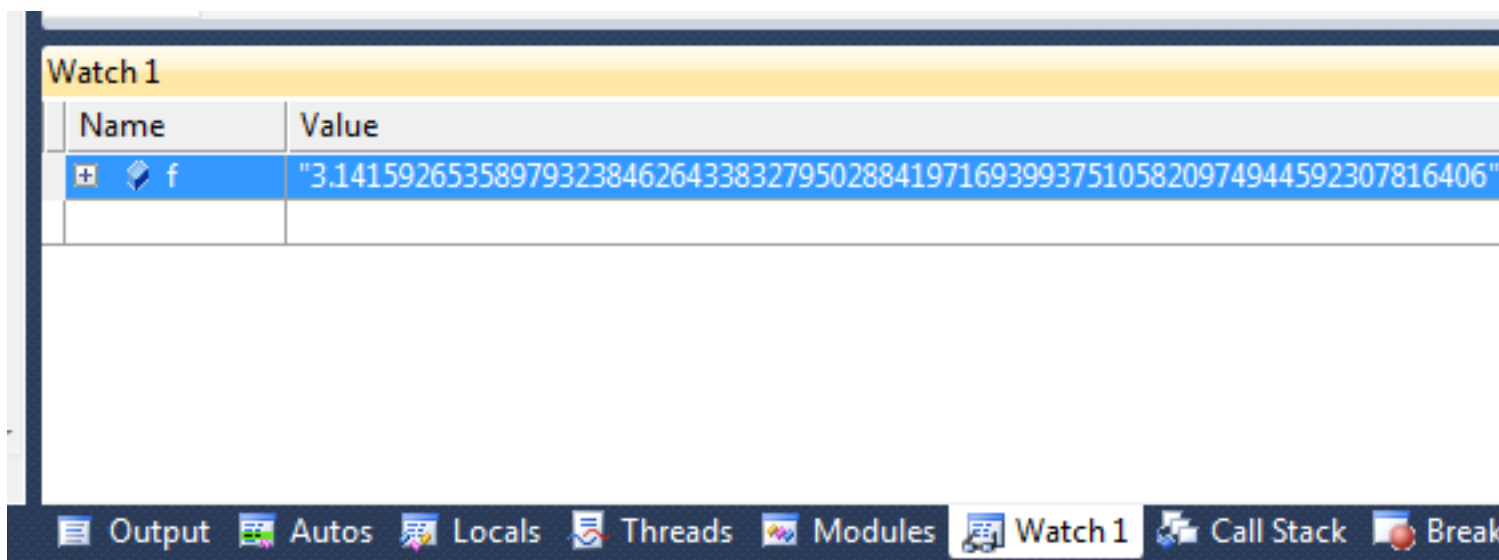


### Note

These visualizers have only been tested with VC10, also given the ability of buggy visualizers to crash your Visual C++ debugger, make sure you back up `autoexp.dat` file before using these!!

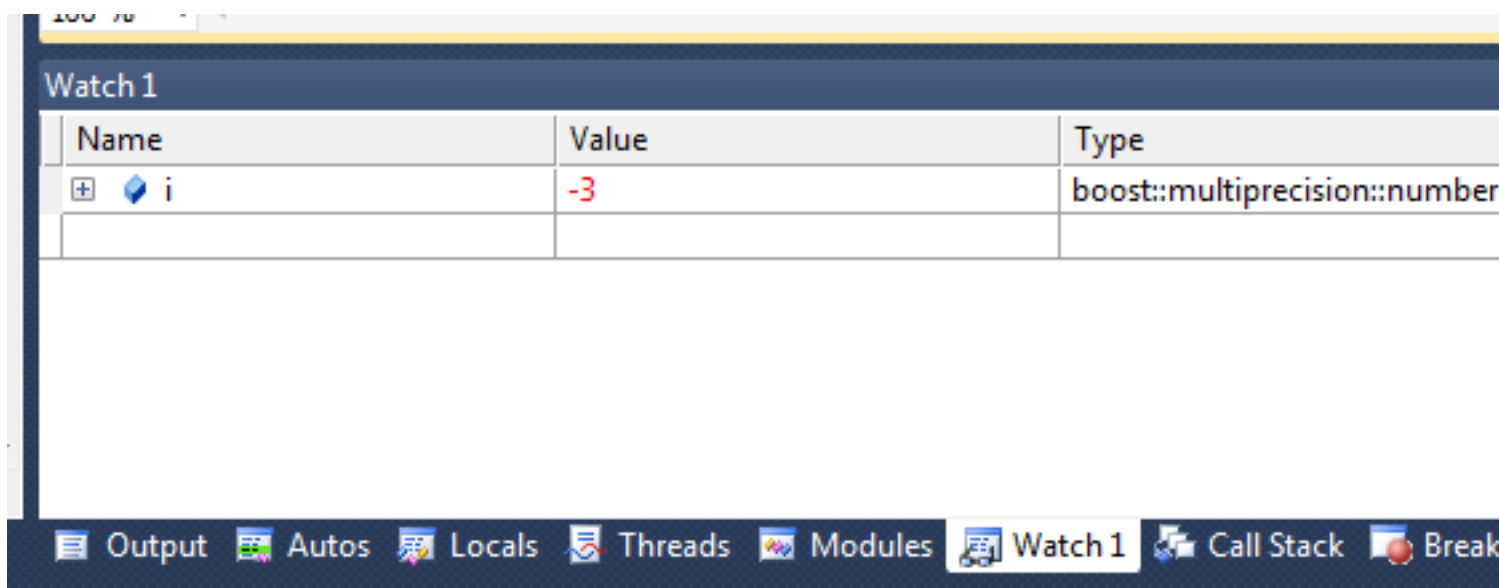
The first visualizer provides improved views of `debug_adaptor`:



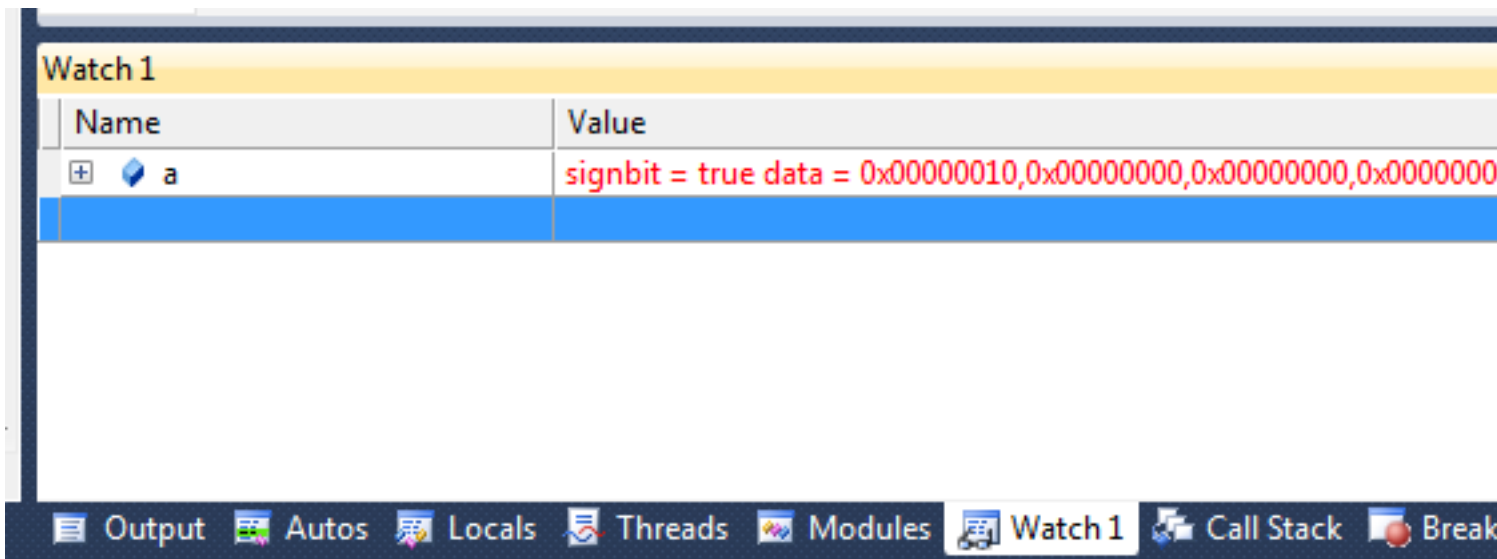


The next visualizer provides improved views of `cpp_int`: small numbers are displayed as actual values, while larger numbers are displayed as an array of hexadecimal parts, with the most significant part first.

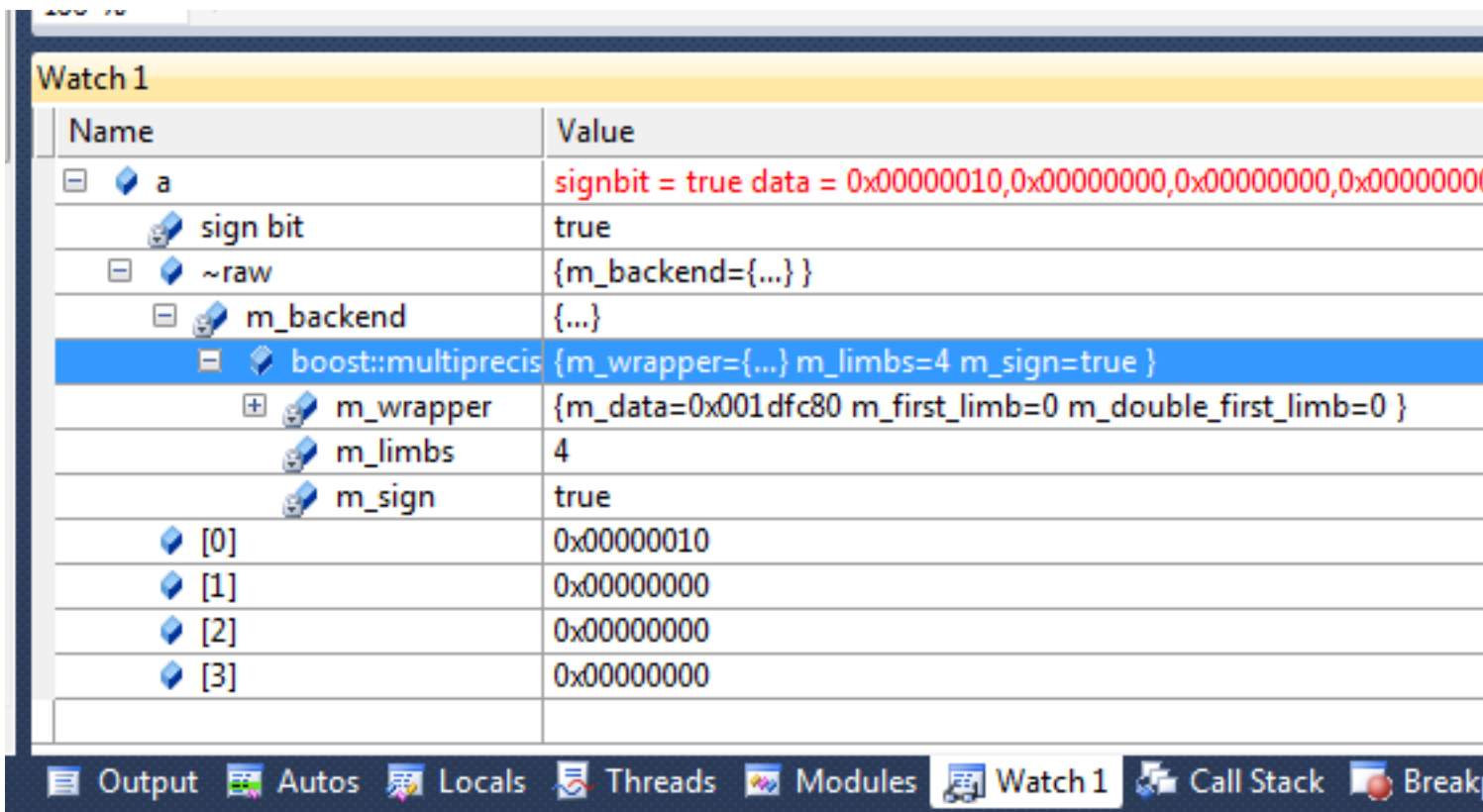
Here's what it looks like for small values:



And for larger values:



There is also a `~raw` child member that lets you see the actual members of the class:



The visualizer for `cpp_dec_float` shows the first few digits of the value in the preview field, and the full array of digits when you expand the view. As before the `~raw` child gives you access to the actual data members:

Watch 1	
Name	Value
dec	-3.1415926535897932...
exponent	0
sign bit	true
+ ~raw	{m_backend={...}}
[0]	3
[1]	14159265
[2]	35897932
[3]	38462643
[4]	38327950
[5]	28841971
[6]	69399375
[7]	10582097
[8]	49445923

## Constructing and Interconverting Between Number Types

All of the number types that are based on number have certain conversion rules in common. In particular:

- Any number type can be constructed (or assigned) from any builtin arithmetic type, as long as the conversion isn't lossy (for example float to int conversion):

```
cpp_dec_float_50 df(0.5);    // OK construction from double
cpp_int          i(450);     // OK constructs from signed int
cpp_int          j = 3.14;   // Error, lossy conversion.
```

- A number can be explicitly constructed from an arithmetic type, even when the conversion is lossy:

```
cpp_int          i(3.14);    // OK explicit conversion
i = static_cast<cpp_int>(3.14) // OK explicit conversion
i.assign(3.14);              // OK, explicit assign and avoid a temporary from the cast above
i = 3.14;                    // Error, no implicit assignment operator for lossy conversion.
cpp_int          j = 3.14;   // Error, no implicit constructor for lossy conversion.
```

- A number can be converted to any built in type, via the `convert_to` member function:

```
mpz_int z(2);
int i = z.template convert_to<int>(); // sets i to 2
```

Additional conversions may be supported by particular backends.

- A number can be converted to any built in type, via an explicit conversion operator: this functionality is only available on compilers supporting C++11's explicit conversion syntax.

```
mpz_int z(2);
int i = z; // Error, implicit conversion not allowed.
int j = static_cast<int>(z); // OK explicit conversion.
```

- Any number type can be *explicitly* constructed (or assigned) from a `const char*` or a `std::string`:

```
// pi to 50 places from a string:
cpp_dec_float_50 df("3.14159265358979323846264338327950288419716939937510");
// Integer type will automatically detect "0x" and "0" prefixes and parse the string accordingly:
cpp_int i("0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF0000000000000000");
// Invalid input always results in a std::runtime_error being thrown:
i = static_cast<cpp_int>("3.14");
// implicit conversions from strings are not allowed:
i = "23"; // Error, no assignment operator for implicit conversion from string
// assign member function, avoids having to create a temporary via a static_cast:
i.assign("23"); // OK
```

- Any number type will interoperate with the builtin types in arithmetic expressions as long as the conversions are not lossy:

```
// pi to 50 places from a string:
cpp_dec_float_50 df = "3.14159265358979323846264338327950288419716939937510";
// Multiply by 2 - using an integer literal here is usually more efficient
// than constructing a temporary:
df *= 2;

// You can't mix integer types with floats though:
cpp_int i = 2;
i *= 3.14; // Error, no *= operator will be found.
```

- Any number type can be streamed to and from the C++ iostreams:

```
cpp_dec_float_50 df = "3.14159265358979323846264338327950288419716939937510";
// Now print at full precision:
std::cout << std::setprecision(std::numeric_limits<cpp_dec_float_50>::max_digits10)
    << df << std::endl;
cpp_int i = 1;
i <= 256;
// Now print in hex format with prefix:
std::cout << std::hex << std::showbase << i << std::endl;
```

- Interconversions between number types of the same family are allowed and are implicit conversions if no loss of precision is involved, and explicit if it is:

```

int128_t      i128 = 0;
int266_t      i256 = i128; // OK implicit widening conversion
int128_t      = i256; // Error, no assignment operator found, narrowing conversion is ex-
plicit
int128_t      = static_cast<int128_t>(i256); // OK, explicit narrowing conversion

mpz_int       z     = 0;
mpf_float     f     = z; // OK, GMP handles this conversion natively, and it's not lossy and
therefore implicit

mpf_float_50 f50    = 2;
f              = f50; // OK, conversion from fixed to variable precision, f will have 50
digits precision.
f50            = f; // Error, conversion from variable to fixed precision is potentially
lossy, explicit cast required.

```

- Some interconversions between number types are completely generic, and are always available, albeit the conversions are always *explicit*:

```

cpp_int cppi(2);
// We can always convert between numbers of the same category -
// int to int, rational to rational, or float to float, so this is OK
// as long as we use an explicit conversion:
mpz_int z(cppi);
// We can always promote from int to rational, int to float, or rational to float:
cpp_rational cpr(cppi); // OK, int to rational
cpp_dec_float_50 df(cppi); // OK, int to float
df              = static_cast<cpp_dec_float_50>(cpr); // OK, explicit rational to float
conversion
// However narrowing and/or implicit conversions always fail:
cppi            = df; // Compiler error, conversion not allowed

```

- Other interconversions may be allowed as special cases, whenever the backend allows it:

```

mpf_t      m; // Native GMP type.
mpf_init_set_ui(m, 0); // set to a value;
mpf_float i(m); // copies the value of the native type.

```

More information on what additional types a backend supports conversions from are given in the tutorial for each backend. The converting constructor will be implicit if the backend's converting constructor is also implicit, and explicit if the backends converting constructor is also explicit.

## Generating Random Numbers

Random numbers are generated in conjunction with Boost.Random. However, since Boost.Random is unaware of [arbitrary precision](#) numbers, it's necessary to include the header:

```
#include <boost/multiprecision/random.hpp>
```

In order to act as a bridge between the two libraries.

Integers with  $N$  random bits are generated using `independent_bits_engine`:

```
#include <boost/multiprecision/gmp.hpp>
#include <boost/multiprecision/random.hpp>

using namespace boost::multiprecision;
using namespace boost::random;

//
// Declare our random number generator type, the underlying generator
// is the Mersenne twister mt19937 engine, and 256 bits are generated:
//
typedef independent_bits_engine<mt19937, 256, mpz_int> generator_type;
generator_type gen;
//
// Generate some values:
//
std::cout << std::hex << std::showbase;
for(unsigned i = 0; i < 10; ++i)
    std::cout << gen() << std::endl;
```

Alternatively we can generate integers in a given range using `uniform_int_distribution`, this will invoke the underlying engine multiple times to build up the required number of bits in the result:

```
#include <boost/multiprecision/gmp.hpp>
#include <boost/multiprecision/random.hpp>

using namespace boost::multiprecision;
using namespace boost::random;

//
// Generate integers in a given range using uniform_int,
// the underlying generator is invoked multiple times
// to generate enough bits:
//
mt19937 mt;
uniform_int_distribution<mpz_int> ui(0, mpz_int(1) << 256);
//
// Generate the numbers:
//
std::cout << std::hex << std::showbase;
for(unsigned i = 0; i < 10; ++i)
    std::cout << ui(mt) << std::endl;
```

Floating point values in  $[0,1)$  are generated using `uniform_01`, the trick here is to ensure that the underlying generator produces as many random bits as there are digits in the floating point type. As above `independent_bits_engine` can be used for this purpose, note that we also have to convert decimal digits (in the floating point type) to bits (in the random number generator):

```
#include <boost/multiprecision/gmp.hpp>
#include <boost/multiprecision/random.hpp>

using namespace boost::multiprecision;
using namespace boost::random;
//
// We need an underlying generator with at least as many bits as the
// floating point type to generate numbers in [0, 1) with all the bits
// in the floating point type randomly filled:
//
uniform_01<mpf_float_50> uf;
independent_bits_engine<mt19937, 50L*1000L/301L, mpz_int> gen;
//
// Generate the values:
//
std::cout << std::setprecision(50);
for(unsigned i = 0; i < 20; ++i)
    std::cout << uf(gen) << std::endl;
```

Finally, we can modify the above example to produce numbers distributed according to some distribution:

```
#include <boost/multiprecision/gmp.hpp>
#include <boost/multiprecision/random.hpp>

using namespace boost::multiprecision;
using namespace boost::random;
//
// We can repeat the above example, with other distributions:
//
uniform_real_distribution<mpf_float_50> ur(-20, 20);
gamma_distribution<mpf_float_50> gd(20);
independent_bits_engine<mt19937, 50L*1000L/301L, mpz_int> gen;
//
// Generate some values:
//
std::cout << std::setprecision(50);
for(unsigned i = 0; i < 20; ++i)
    std::cout << ur(gen) << std::endl;
for(unsigned i = 0; i < 20; ++i)
    std::cout << gd(gen) << std::endl;
```

## Primality Testing

The library implements a Miller-Rabin test for primality:

```
#include <boost/multiprecision/miller_rabin.hpp>

template <class Backend, expression_template_option ExpressionTemplates, class Engine>
bool miller_rabin_test(const number<Backend, ExpressionTemplates>& n, unsigned trials, Engine& gen);

template <class Backend, expression_template_option ExpressionTemplates, class Engine>
bool miller_rabin_test(const number<Backend, ExpressionTemplates>& n, unsigned trials);
```

These functions perform a Miller-Rabin test for primality, if the result is `false` then  $n$  is definitely composite, while if the result is `true` then  $n$  is probably prime. The probability to declare a composite  $n$  as probable prime is at most  $0.25^{\text{trials}}$ . Note that this does not allow a statement about the probability of  $n$  being actually prime (for that, the prior probability would have to be known). The algorithm used performs some trial divisions to exclude small prime factors, does one Fermat test to exclude many more composites, and then

uses the Miller-Rabin algorithm straight out of Knuth Vol 2, which recommends 25 trials for a pretty strong likelihood that  $n$  is prime.

The third optional argument is for a Uniform Random Number Generator from Boost.Random. When not provided the mt19937 generator is used. Note that when producing random primes then you should probably use a different random number generator to produce candidate prime numbers for testing, than is used internally by `miller_rabin_test` for determining whether the value is prime. It also helps of course to seed the generators with some source of randomness.

The following example searches for a prime  $p$  for which  $(p-1)/2$  is also probably prime:

```
#include <boost/multiprecision/cpp_int.hpp>
#include <boost/multiprecision/miller_rabin.hpp>
#include <iostream>
#include <iomanip>

int main()
{
    using namespace boost::random;
    using namespace boost::multiprecision;

    typedef cpp_int int_type;
    mt11213b base_gen(clock());
    independent_bits_engine<mt11213b, 256, int_type> gen(base_gen);
    //
    // We must use a different generator for the tests and number generation, otherwise
    // we get false positives.
    //
    mt19937 gen2(clock());

    for(unsigned i = 0; i < 100000; ++i)
    {
        int_type n = gen();
        if(miller_rabin_test(n, 25, gen2))
        {
            // Value n is probably prime, see if (n-1)/2 is also prime:
            std::cout << "We have a probable prime with value: " << std::hex << std::showbase << n << std::endl;
            if(miller_rabin_test((n-1)/2, 25, gen2))
            {
                std::cout << "We have a safe prime with value: " << std::hex << std::showbase << n << std::endl;
                return 0;
            }
        }
    }
    std::cout << "Oops, no safe primes were found" << std::endl;
    return 1;
}
```

## Literal Types and `constexpr` Support



### Note

The features described in this section make heavy use of C++11 language features, currently (as of May 2013) only GCC-4.7 and later, and Clang 3.3 and later have the support required to make these features work.

There is limited support for `constexpr` and user-defined literals in the library, currently the number front end supports `constexpr` on default construction and all forwarding constructors, but not on any of the non-member operators. So if some type  $B$  is a literal



type, then `number<B>` is also a literal type, and you will be able to compile-time-construct such a type from any literal that `B` is compile-time-constructible from. However, you will not be able to perform compile-time arithmetic on such types.

Currently the only backend type provided by the library that is also a literal type are instantiations of `cpp_int_backend` where the Allocator parameter is type `void`, and the Checked parameter is `boost::multiprecision::unchecked`.

For example:

```
using namespace boost::multiprecision;

constexpr int128_t      i = 0;      // OK, fixed precision int128_t has no allocator.
constexpr uint1024_t   j = 0xFFFFFFFF00000000uLL; // OK, fixed precision uint1024_t has no allocator.

constexpr checked_uint128_t k = -1; // Error, checked type is not a literal type as we need runtime error checking.
constexpr cpp_int         l = 2;    // Error, type is not a literal as it performs memory management.
```

There is also limited support for user defined-literals - these are limited to unchecked, fixed precision `cpp_int`'s which are specified in hexadecimal notation. The suffixes supported are:

Suffix	Meaning
<code>_cppi</code>	Specifies a value of type: <code>number&lt;cpp_int_backend&lt;N,N,signed_magnitude,unchecked,void&gt;&gt;</code> , where <code>N</code> is chosen to contain just enough digits to hold the number specified.
<code>_cppui</code>	Specifies a value of type: <code>number&lt;cpp_int_backend&lt;N,N,unsigned_magnitude,unchecked,void&gt;&gt;</code> , where <code>N</code> is chosen to contain just enough digits to hold the number specified.
<code>_cppiN</code>	Specifies a value of type <code>number&lt;cpp_int_backend&lt;N,N,signed_magnitude,unchecked,void&gt;&gt;</code> .
<code>_cppuiN</code>	Specifies a value of type <code>number&lt;cpp_int_backend&lt;N,N,unsigned_magnitude,unchecked,void&gt;&gt;</code> .

In each case, use of these suffixes with hexadecimal values produces a `constexpr` result.

Examples:

```
//
// Any use of user defined literals requires that we import the literal-operators
// into current scope first:
using namespace boost::multiprecision::literals;
//
// To keep things simple in the example, we'll make our types used visible to this scope as well:
using namespace boost::multiprecision;
//
// The value zero as a number<cpp_int_backend<4,4,signed_magnitude,unchecked,void>> >:
constexpr auto a = 0x0_cppi;
// The type of each constant has 4 bits per hexadecimal digit,
// so this is of type uint256_t (ie number<cpp_int_backend<256,256,unsigned_magnitude,unJ
checked,void>>):
constexpr auto b = 0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF_cppui;
//
// Smaller values can be assigned to larger values:
int256_t c = 0x1234_cppi; // OK
//
// However, this does not currently work in constexpr contexts:
constexpr int256_t d = 0x1_cppi; // Compiler error
//
// Constants can be padded out with leading zeros to generate wider types:
constexpr uint256_t e = 0x0000000000000000000000000000000000000000000000000000000000000000_FFFFFFFF_cpl
pui; // OK
//
// However, specific width types are best produced with specific-width suffixes,
// ones supported by default are `_cpp[ui]i128`, `_cpp[ui]i256`, `_cpp[ui]i512`, `_cpp[ui]i1024`.
//
constexpr int128_t f = 0x1234_cppii128; // OK, always produces an int128_t as the result.
constexJ
pr uint1024_t g = 0xaaabbbbbbbcccccccccccccccccccccccccccccccccccccccccccccccccccccccccccccccc_cpJ
puil1024;
//
// If other specific width types are required, then there is a macro for generating the operators
// for these. The macro can be used at namespace scope only:
//
BOOST_MP_DEFINE_SIZED_CPP_INT_LITERAL(2048);
//
// Now we can create 2048-bit literals as well:
constexpr auto h = 0xff_cppii2048; // h is of type number<cpp_int_backend<2048,2048,signed_mag,J
nitude,unchecked,void>> >
//
// Finally negative values are handled via the unary minus operator:
//
constexpr int1024_t i = -0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF_cpl
puil1024;
//
// Which means this also works:
constexpr int1024_t j = -g; // OK: unary minus operator is constexpr.
```

## Rounding Rules for Conversions

As a general rule, all conversions between unrelated types are performed using basic arithmetic operations, therefore conversions are either exact, or follow the same rounding rules as arithmetic for the type in question.

The following table summarises the situation for conversions from native types:

Backend	Rounding Rules
<a href="#">cpp_int</a>	Conversions from integer types are exact if the target has sufficient precision, otherwise they truncate to the first $2^{\text{MaxBits}}$ bits (modulo arithmetic). Conversions from floating point types are truncating to the nearest integer.
<a href="#">gmp_int</a>	Conversions are performed by the GMP library except for conversion from <code>long double</code> which is truncating.
<a href="#">tom_int</a>	Conversions from floating point types are truncating, all others are performed by libtommath and are exact.
<a href="#">gmp_float</a>	Conversions are performed by the GMP library except for conversion from <code>long double</code> which should be exact provided the target type has as much precision as a <code>long double</code> .
<a href="#">mpfr_float</a>	All conversions are performed by the underlying MPFR library.
<a href="#">cpp_dec_float</a>	All conversions are performed using basic arithmetic operations and are truncating.
<a href="#">gmp_rational</a>	See <a href="#">gmp_int</a>
<a href="#">cpp_rational</a>	See <a href="#">cpp_int</a>
<a href="#">tommath_rational</a>	See <a href="#">tom_int</a>

## Mixed Precision Arithmetic

Mixed precision arithmetic is fully supported by the library.

There are two different forms:

- Where the operands are of different precision.
- Where the operands are of the same precision, but yield a higher precision result.

### Mixing Operands of Differing Precision

If the arguments to a binary operator are of different precision, then the operation is allowed as long as there is an unambiguous implicit conversion from one argument type to the other. In all cases the arithmetic is performed "as if" the lower precision type is promoted to the higher precision type before applying the operator. However, particular backends may optimise this and avoid actually creating a temporary if they are able to do so.

For example:

```
mpfr_float_50      a(2), b;
mpfr_float_100     c(3), d;
static_mpfr_float_50 e(5), f;
mpz_int            i(20);

d = a * c; // OK, result of operand is an mpfr_float_100.
b = a * c; // Error, can't convert the result to an mpfr_float_50 as it will lose digits.
f = a * e;  // Error, operator is ambiguous, result could be of either type.
f = e * i;  // OK, unambiguous conversion from mpz_int to static_mpfr_float_50
```

## Operands of the Same Precision

Sometimes you want to apply an operator to two arguments of the same precision in such a way as to obtain a result of higher precision. The most common situation occurs with fixed precision integers, where you want to multiply two  $N$ -bit numbers to obtain a  $2N$ -bit result. This is supported in this library by the following free functions:

```
template <class ResultType, class Source1 class Source2>
ResultType& add(ResultType& result, const Source1& a, const Source2& b);

template <class ResultType, class Source1 class Source2>
ResultType& subtract(ResultType& result, const Source1& a, const Source2& b);

template <class ResultType, class Source1 class Source2>
ResultType& multiply(ResultType& result, const Source1& a, const Source2& b);
```

These functions apply the named operator to the arguments *a* and *b* and store the result in *result*, returning *result*. In all cases they behave "as if" arguments *a* and *b* were first promoted to type `ResultType` before applying the operator, though particular backends may well avoid that step by way of an optimization.

The type `ResultType` must be an instance of class `number`, and the types `Source1` and `Source2` may be either instances of class `number` or native integer types. The latter is an optimization that allows arithmetic to be performed on native integer types producing an extended precision result.

For example:

```
#include <boost/multiprecision/cpp_int.hpp>

using namespace boost::multiprecision;

boost::uint64_t i = (std::numeric_limits<boost::uint64_t>::max)();
boost::uint64_t j = 1;

uint128_t ui128;
uint256_t ui256;
//
// Start by performing arithmetic on 64-bit integers to yield 128-bit results:
//
std::cout << std::hex << std::showbase << i << std::endl;
std::cout << std::hex << std::showbase << add(ui128, i, j) << std::endl;
std::cout << std::hex << std::showbase << multiply(ui128, i, i) << std::endl;
//
// The try squaring a 128-bit integer to yield a 256-bit result:
//
ui128 = (std::numeric_limits<uint128_t>::max)();
std::cout << std::hex << std::showbase << multiply(ui256, ui128, ui128) << std::endl;
```

Produces the output:

```
0xffffffffffffffff
0x10000000000000000
0xFFFFFFFFFFFFFFFF0000000000000001
0xFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF000000000000000000000001
```

## Backends With Optimized Mixed Precision Arithmetic

The following backends have at least some direct support for mixed precision arithmetic, and therefore avoid creating unnecessary temporaries when using the interfaces above. Therefore when using these types it's more efficient to use mixed precision arithmetic, than it is to explicitly cast the operands to the result type:

[mpfr\\_float](#), [gmp\\_float](#), [cpp\\_int](#).

## Generic Integer Operations

All of the [non-member integer operations](#) are overloaded for the built in integer types in `<boost/multiprecision/integer.hpp>`. Where these operations require a temporary increase in precision (such as for `powm`), then if no built in type is available, a [cpp\\_int](#) of appropriate precision will be used.

Some of these functions are trivial, others use compiler intrinsics (where available) to ensure optimal evaluation.

The overloaded functions are:

```
template <class Integer, class I2>
Integer& multiply(Integer& result, const I2& a, const I2& b);
```

Multiplies two I2 values, to produce a wider Integer result.

Returns `result = a * b` without overflow or loss of precision in the multiplication.

```
template <class Integer, class I2>
Integer& add(Integer& result, const I2& a, const I2& b);
```

Adds two I2 values, to produce a wider Integer result.

Returns `result = a + b` without overflow or loss of precision in the addition.

```
template <class Integer, class I2>
Integer& subtract(Integer& result, const I2& a, const I2& b);
```

Subtracts two I2 values, to produce a wider Integer result.

Returns `result = a - b` without overflow or loss of precision in the subtraction.

```
template <class Integer>
Integer powm(const Integer& b, const Integer& p, const Integer& m);
```

Returns  $b^p \% m$ .

```
template <class Integer>
void divide_qr(const Integer& x, const Integer& y, Integer& q, Integer& r);
```

Sets `q = x / y` and `r = x % y`.

```
template <class Integer1, class Integer2>
Integer2 integer_modulus(const Integer1& x, Integer2 val);
```

Returns `x % val`;

```
template <class Integer>
unsigned lsb(const Integer& x);
```

Returns the (zero-based) index of the least significant bit of `x`.

Throws a `std::domain_error` if `x <= 0`.

```
template <class Integer>
unsigned msb(const Integer& x);
```

Returns the (zero-based) index of the most significant bit of  $x$ .

Throws a `std::domain_error` if  $x \leq 0$ .

```
template <class Integer>
bool bit_test(const Integer& val, unsigned index);
```

Returns true if bit index is set in  $val$ .

```
template <class Integer>
Integer& bit_set(Integer& val, unsigned index);
```

Sets the index bit in  $val$ .

```
template <class Integer>
Integer& bit_unset(Integer& val, unsigned index);
```

Unsets the index bit in  $val$ .

```
template <class Integer>
Integer& bit_flip(Integer& val, unsigned index);
```

Flips the index bit in  $val$ .

```
template <class Integer>
Integer sqrt(const Integer& x);
template <class Integer>
Integer sqrt(const Integer& x, Integer& r);
```

Returns the integer square root  $s$  of  $x$  and sets  $r$  to the remainder  $x - s^2$ .

```
template <class Engine>
bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials, Engine& gen);
bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials);
```

The regular Miller-Rabin functions in `<boost/multiprecision/miller_rabin.hpp>` are defined in terms of the above generic operations, and so function equally well for built in and multiprecision types.

## Boost.Serialization Support

Support for serialization comes in two forms:

- Classes `number`, `debug_adaptor`, `logged_adaptor` and `rational_adaptor` have "pass through" serialization support which requires the underlying backend to be serializable.
- Backends `cpp_int`, `cpp_bin_float`, `cpp_dec_float` and `float128` have full support for Boost.Serialization.

## Numeric Limits

Boost.Multiprecision tries hard to implement `std::numeric_limits` for all types as far as possible and meaningful because experience with Boost.Math has shown that this aids portability.

The [C++ standard library](#) defines `std::numeric_limits` in section 18.3.2.

This in turn refers to the C standard [SC22/WG11 N507 DRAFT INTERNATIONAL ISO/IEC STANDARD WD 10967-1](#) Information technology Language independent arithmetic Part 1: Integer and Floating point arithmetic.

That C Standard in turn refers to

[IEEE754 IEEE Standard for Binary Floating-Point Arithmetic](#)

There is a useful summary at [C++ reference](#).

The chosen backend often determines how completely `std::numeric_limits` is available.

Compiler options, processor type, and definition of macros or assembler instructions to control denormal numbers will alter the values in the tables given below.



### Warning

GMP's `mpf_t` does not have a concept of overflow: operations that lead to overflow eventually run out of resources and terminate with stack overflow (often after several seconds).

## `std::numeric_limits<> constants`

### `is_specialized`

true for all arithmetic types (integer, floating and fixed-point) for which `std::numeric_limits<T>::numeric_limits` is specialized.

A typical test is

```
if (std::numeric_limits<T>::is_specialized == false)
{
    std::cout << "type " << typeid(T).name() << " is not specialized for std::numeric_lim
its!" << std::endl;
    // ...
}
```

Typically `numeric_limits<T>::is_specialized` is true for all `T` where the compile-time constant members of `numeric_limits` are indeed known at compile time, and don't vary at runtime. For example floating point types with runtime-variable precision such as `mpfr_float` have no `numeric_limits` specialization as it would be impossible to define all the members at compile time. In contrast the precision of a type such as `mpfr_float_50` is known at compile time, and so it *does* have a `numeric_limits` specialization.

Note that not all the `std::numeric_limits` member constants and functions are meaningful for all user-defined types (UDT), such as the decimal and binary multiprecision types provided here. More information on this is given in the sections below.

### `infinity`

For floating-point types,  $\infty$  is defined wherever possible, but clearly infinity is meaningless for `__arbitrary_precision` arithmetic backends, and there is one floating point type (GMP's `mpf_t`, see [gmp\\_float](#)) which has no notion of infinity or NaN at all.

A typical test whether infinity is implemented is

```
if (std::numeric_limits<T>::has_infinity)
{
    std::cout << std::numeric_limits<T>::infinity() << std::endl;
}
```

and using tests like this is strongly recommended to improve portability.

If the backend is switched to a type that does not support infinity then, without checks like this, there will be trouble.

## is\_signed

`std::numeric_limits<T>::is_signed == true` if the type `T` is signed.

For built-in binary types, the sign is held in a single bit, but for other types (`cpp_dec_float` and `cpp_bin_float`) it may be a separate storage element, usually `bool`.

## is\_exact

`std::numeric_limits<T>::is_exact == true` if type `T` uses exact representations.

This is defined as `true` for all integer types and `false` for floating-point types.

[A usable definition](#) has been discussed.

ISO/IEC 10967-1, Language independent arithmetic, noted by the C++ Standard defines

A floating point type `F` shall be a finite subset of `[real]`.

The important practical distinction is that all integers (up to `max()`) can be stored exactly.

[Rational](#) types using two integer types are also exact.

Floating-point types **cannot store all real values** (those in the set of  $\mathbb{R}$ ) **exactly**. For example, 0.5 can be stored exactly in a binary floating-point, but 0.1 cannot. What is stored is the nearest representable real value, that is, rounded to nearest.

Fixed-point types (usually decimal) are also defined as exact, in that they only store a **fixed precision**, so half cents or pennies (or less) cannot be stored. The results of computations are rounded up or down, just like the result of integer division stored as an integer result.

There are number of proposals to [add Decimal Floating Point Support to C++](#).

[Decimal TR](#).

And also [C++ Binary Fixed-Point Arithmetic](#).

## is\_bounded

`std::numeric_limits<T>::is_bounded == true` if the set of values represented by the type `T` is finite.

This is `true` for all built-in integer, fixed and floating-point types, and most multi-precision types.

It is only `false` for a few `__arbitrary_precision` types like `cpp_int`.

Rational and fixed-exponent representations are exact but not integer.

## is\_modulo

`std::numeric_limits<T>::is_modulo` is defined as `true` if adding two positive values of type `T` can yield a result less than either value.

`is_modulo == true` means that the type does not overflow, but, for example, 'wraps around' to zero, when adding one to the `max()` value.

For most built-in integer types, `std::numeric_limits<>::is_modulo` is `true`.

`bool` is the only exception.



The modulo behaviour is sometimes useful, but also can be unexpected, and sometimes undesired, behaviour.

Overflow of signed integers can be especially unexpected, possibly causing change of sign.

Boost.Multiprecision integer type `cpp_int` is not modulo because as an `__arbitrary_precision` types, it expands to hold any value that the machine resources permit.

However fixed precision `cpp_int`'s may be modulo if they are unchecked (i.e. they behave just like built in integers), but not if they are checked (overflow causes an exception to be raised).

Built-in and multi-precision floating-point types are normally not modulo.

Where possible, overflow is to `std::numeric_limits<>::infinity()`, provided `std::numeric_limits<>::has_infinity == true`.

## radix

Constant `std::numeric_limits<T>::radix` returns either 2 (for built-in and binary types) or 10 (for decimal types).

## digits

The number of `radix` digits that be represented without change:

- for integer types, the number of **non-sign bits** in the significand.
- for floating types, the number of **radix digits** in the significand.

The values include any implicit bit, so for example, for the ubiquitous `double` using 64 bits ([IEEE binary64](#)), `digits == 53`, even though there are only 52 actual bits of the significand stored in the representation. The value of `digits` reflects the fact that there is one implicit bit which is always set to 1.

The Boost.Multiprecision binary types do not use an implicit bit, so the `digits` member reflects exactly how many bits of precision were requested:

```
typedef number<cpp_bin_float<53, digit_base_2> > float64;
typedef number<cpp_bin_float<113, digit_base_2> > float128;
std::numeric_limits<float64>::digits == 53.
std::numeric_limits<float128>::digits == 113.
```

For the most common case of `radix == 2`, `std::numeric_limits<T>::digits` is the number of bits in the representation, not counting any sign bit.

For a decimal integer type, when `radix == 10`, it is the number of decimal digits.

## digits10

Constant `std::numeric_limits<T>::digits10` returns the number of decimal digits that can be represented without change or loss.

For example, `numeric_limits<unsigned char>::digits10` is 2.

This somewhat inscrutable definition means that an `unsigned char` can hold decimal values 0..99 without loss of precision or accuracy, usually from truncation.

Had the definition been 3 then that would imply it could hold 0..999, but as we all know, an 8-bit `unsigned char` can only hold 0..255, and an attempt to store 256 or more will involve loss or change.

For bounded integers, it is thus **one less** than number of decimal digits you need to display the biggest integer `std::numeric_limits<T>::max()`. This value can be used to predict the layout width required for

```
std::cout
<< std::setw(std::numeric_limits<short>::digits10 +1 +1) // digits10+1, and +1 for sign.
<< std::showpos << (std::numeric_limits<short>::max)() // +32767
<< std::endl
<< std::setw(std::numeric_limits<short>::digits10 +1 +1)
<< (std::numeric_limits<short>::min)() << std::endl; // -32767
```

For example, unsigned short is often stored in 16 bits, so the maximum value is 0xFFFF or 65535.

```
std::cout
<< std::setw(std::numeric_limits<unsigned short>::digits10 +1 +1) // digits10+1, and +1 for sign.
<< std::showpos << (std::numeric_limits<unsigned short>::max)() // 65535
<< std::endl
<< std::setw(std::numeric_limits<unsigned short>::digits10 +1 +1) // digits10+1, and +1 for sign.
<< (std::numeric_limits<unsigned short>::min)() << std::endl; // 0
```

For bounded floating-point types, if we create a double with a value with digits10 (usually 15) decimal digits, 1e15 or 1000000000000000 :

```
std::cout.precision(std::numeric_limits<double>::max_digits10);
double d = 1e15;
double dpl = d+1;
std::cout << d << "\n" << dpl << std::endl;
// 1000000000000000
// 10000000000000001
std::cout << dpl - d << std::endl; // 1
```

and we can increment this value to 10000000000000001 as expected and show the difference too.

But if we try to repeat this with more than digits10 digits,

```
std::cout.precision(std::numeric_limits<double>::max_digits10);
double d = 1e16;
double dpl = d+1;
std::cout << d << "\n" << dpl << std::endl;
// 10000000000000000
// 10000000000000000
std::cout << dpl - d << std::endl; // 0 !!!
```

then we find that when we add one it has no effect, and display show that there is loss of precision. See [Loss of significance or cancellation error](#).

So digits10 is the number of decimal digits **guaranteed** to be correct.

For example, 'round-tripping' for double:

- If a decimal string with at most digits10( == 15) significant decimal digits is converted to double and then converted back to the same number of significant decimal digits, then the final string will match the original 15 decimal digit string.
- If a double floating-point number is converted to a decimal string with at least 17 decimal digits and then converted back to double, then the result will be binary identical to the original double value.

For most purposes, you will much more likely want std::numeric\_limits<>::max\_digits10, the number of decimal digits that ensure that a change of one least significant bit (ULP) produces a different decimal digits string.

For nearly all floating-point types, max\_digits10 is digits10+2, but you should use max\_digits10 where possible.

If max\_digits10 is not available, you should using the [Kahan formula for floating-point type T](#)

```
max_digits10 = std::numeric_limits<T>::digits * 3010U/10000U;
```

The factor is  $\log_{10}(2) = 0.3010$  but must be evaluated at compile time using only integers.

(See also [Richard P. Brent and Paul Zimmerman, Modern Computer Arithmetic](#) Equation 3.8 on page 116.).

The extra two (or 3) least significant digits are 'noisy' and may be junk, but if you want to 'round-trip' - printing a value out and reading it back in - you must use `os.precision(std::numeric_limits<T>::max_digits10)`. For at least one popular compiler, you must also use `std::scientific` format.

## max\_digits10

`std::numeric_limits<T>::max_digits10` was added for floating-point because `digits10` decimal digits are insufficient to show a least significant bit (ULP) change giving puzzling displays like

```
0.6666666666666667 != 0.6666666666666667
```

from failure to 'round-trip', for example:

```
double write = 2./3; // Any arbitrary value that cannot be represented exactly.
double read = 0;
std::stringstream s;
s.precision(std::numeric_limits<double>::digits10); // or `float64_t` for 64-bit IEEE754 double.
s << write;
s >> read;
if(read != write)
{
    std::cout << std::setprecision(std::numeric_limits<double>::digits10)
              << read << " != " << write << std::endl;
}
```

If you wish to ensure that a change of one least significant bit (ULP) produces a different decimal digits string, then `max_digits10` is the precision to use.

For example:

```
double pi = boost::math::double_constants::pi;
std::cout.precision(std::numeric_limits<double>::max_digits10);
std::cout << pi << std::endl; // 3.1415926535897931
```

will display  $\pi$  to the maximum possible precision using a double.

and similarly for a much higher precision type:

```
using namespace boost::multiprecision;

typedef number<cpp_dec_float<50> > cpp_dec_float_50; // 50 decimal digits.

using boost::multiprecision::cpp_dec_float_50;

cpp_dec_float_50 pi = boost::math::constants::pi<cpp_dec_float_50>();
std::cout.precision(std::numeric_limits<cpp_dec_float_50>::max_digits10);
std::cout << pi << std::endl;
// 3.141592653589793238462643383279502884197169399375105820974944592307816406
```

For integer types, `max_digits10` is implementation-dependant, but is usually `digits10 + 2`. This is the output field width required for the maximum value of the type `T` `std::numeric_limits<T>::max()` including a sign and a space.

So this will produce neat columns.

```
std::cout << std::setw(std::numeric_limits<int>::max_digits10) ...
```



### Note

For Microsoft Visual Studio 2010, `std::numeric_limits<float>::max_digits10` is wrongly defined as 8. It should be 9.



### Note

For Microsoft Visual Studio, and default float format, a small range of values approximately 0.0001 to 0.004, with exponent values of 3f2 to 3f6, are wrongly input by one least significant bit, probably every third value of significand.

A workaround is using scientific or exponential format `<< std::scientific`.



### Note

`BOOST_NO_CXX11_NUMERIC_LIMITS` is a suitable feature-test macro to determine if `std::numeric_limits<float>::max_digits10` is implemented on any platform. If `max_digits10` is not available, you should using the [Kahan formula for floating-point type T](#). See above.

For example, to be portable, including older platforms:

```
typedef float T; // Any type: `double`, cpp_dec_float_50, bin_128bit_double_type ...

#if defined(BOOST_NO_CXX11_NUMERIC_LIMITS)
    std::cout.precision(2 + std::numeric_limits<T>::digits * 3010U/10000U);
#else
#   if (_MSC_VER <= 1600) // Correct wrong value for float.
        std::cout.precision(2 + std::numeric_limits<T>::digits * 3010U/10000U);
#   else
        std::cout.precision(std::numeric_limits<T>::max_digits10);
#   endif
#endif

std::cout << "std::cout.precision = " << std::cout.precision() << std::endl;

double x = 1.234567890123456789;

std::cout << "x = " << x << std::endl; //
```

which should output:

```
std::cout.precision = 9
x = 1.23456789
```

## round\_style

The rounding style determines how the result of floating-point operations is treated when the result cannot be **exactly represented** in the significand. Various rounding modes may be provided:

- round to nearest up or down (default for floating-point types).
- round up (toward positive infinity).

- round down (toward negative infinity).
- round toward zero (integer types).
- no rounding (if decimal radix).
- rounding mode is not determinable.

For integer types, `std::numeric_limits<T>::round_style` is always towards zero, so

```
std::numeric_limits<T>::round_style == std::round_to_zero;
```

A decimal type, `cpp_dec_float` rounds in no particular direction, which is to say it doesn't round at all. And since there are several guard digits, it's not really the same as truncation (round toward zero) either.

For floating-point types, it is normal to round to nearest.

```
std::numeric_limits<T>::round_style == std::round_to_nearest;
```

See function `std::numeric_limits<T>::round_error` for the maximum error (in ULP) that rounding can cause.

## has\_denorm\_loss

true if a loss of precision is detected as a [denormalization](#) loss, rather than an inexact result.

Always false for integer types.

false for all types which do not have `has_denorm == std::denorm_present`.

## denorm\_style

[Denormalized values](#) are representations with a variable number of exponent bits that can permit gradual underflow, so that, if type `T` is double.

```
std::numeric_limits<T>::denorm_min() < std::numeric_limits<T>::min()
```

A type may have any of the following enum `float_denorm_style` values:

- `std::denorm_absent`, if it does not allow denormalized values. (Always used for all integer and exact types).
- `std::denorm_present`, if the floating-point type allows denormalized values.
- `std::denorm_indeterminate`, if indeterminate at compile time.

## Tinyness before rounding

```
bool std::numeric_limits<T>::tinyness_before
```

true if a type can determine that a value is too small to be represent as a normalized value before rounding it.

Generally true for `is_iec559` floating-point built-in types, but false for integer types.

Standard-compliant IEEE 754 floating-point implementations may detect the floating-point underflow at three predefined moments:

1. After computation of a result with absolute value smaller than `std::numeric_limits<T>::min()`, such implementation detects *tinyness before rounding* (e.g. UltraSparc).
2. After rounding of the result to `std::numeric_limits<T>::digits` bits, if the result is tiny, such implementation detects *tinyness after rounding* (e.g. SuperSparc).

3. If the conversion of the rounded tiny result to subnormal form resulted in the loss of precision, such implementation detects *denorm loss*.

## std::numeric\_limits<> functions

### max function

Function `std::numeric_limits<T>::max()` returns the largest finite value that can be represented by the type `T`. If there is no such value (and `numeric_limits<T>::bounded` is `false`) then returns `T()`.

For built-in types there is usually a corresponding MACRO value `TYPE_MAX`, where `TYPE` is `CHAR`, `INT`, `FLOAT` etc.

Other types, including those provided by a typedef, for example `INT64_T_MAX` for `int64_t`, may provide a macro definition.

To cater for situations where no `numeric_limits` specialization is available (for example because the precision of the type varies at runtime), packaged versions of this (and other functions) are provided using

```
#include <boost/math/tools/precision.hpp>

T = boost::math::tools::max_value<T>();
```

Of course, these simply use `std::numeric_limits<T>::max()` if available, but otherwise 'do something sensible'.

### lowest function

Since C++11: `std::numeric_limits<T>::lowest()` is

- For integral types, the same as function `min()`.
- For floating-point types, generally the negative of `max()` (but implementation-dependent).

```
-(std::numeric_limits<double>::max)() == std::numeric_limits<double>::lowest();
```

### min function

Function `std::numeric_limits<T>::min()` returns the minimum finite value that can be represented by the type `T`.

For built-in types there is usually a corresponding MACRO value `TYPE_MIN`, where `TYPE` is `CHAR`, `INT`, `FLOAT` etc.

Other types, including those provided by a typedef, for example `INT64_T_MIN` for `int64_t`, may provide a macro definition.

For floating-point types, it is more fully defined as the *minimum positive normalized value*.

See `std::numeric_limits<T>::denorm_min()` for the smallest denormalized value, provided

```
std::numeric_limits<T>::has_denorm == std::denorm_present
```

To cater for situations where no `numeric_limits` specialization is available (for example because the precision of the type varies at runtime), packaged versions of this (and other functions) are provided using

```
#include <boost/math/tools/precision.hpp>

T = boost::math::tools::min_value<T>();
```

Of course, these simply use `std::numeric_limits<T>::min()` if available.

## denorm\_min function

Function `std::numeric_limits<T>::denorm_min()` returns the smallest [denormalized value](#), provided

```
std::numeric_limits<T>::has_denorm == std::denorm_present
```

```
std::cout.precision(std::numeric_limits<double>::max_digits10);
if (std::numeric_limits<double>::has_denorm == std::denorm_present)
{
    double d = std::numeric_limits<double>::denorm_min();

    std::cout << d << std::endl; // 4.9406564584124654e-324

    int exponent;

    double significand = frexp(d, &exponent);
    std::cout << "exponent = " << std::hex << exponent << std::endl; // fffffbcbf
    std::cout << "significand = " << std::hex << significand << std::endl; // 0.500000000000000000
}
else
{
    std::cout << "No denormalization. " << std::endl;
}
```

The exponent is effectively reduced from -308 to -324 (though it remains encoded as zero and leading zeros appear in the significand, thereby losing precision until the significand reaches zero).

## round\_error

Function `std::numeric_limits<T>::round_error()` returns the maximum error (in units of [ULP](#)) that can be caused by any basic arithmetic operation.

```
round_style == std::round_indeterminate;
```

The rounding style is indeterminable at compile time.

For floating-point types, when rounding is to nearest, only half a bit is lost by rounding, and `round_error == 0.5`. In contrast when rounding is towards zero, or plus/minus infinity, we can loose up to one bit from rounding, and `round_error == 1`.

For integer types, rounding always to zero, so at worst almost one bit can be rounded, so `round_error == 1`.

`round_error()` can be used with `std::numeric_limits<T>::epsilon()` to estimate the maximum potential error caused by rounding. For typical floating-point types, `round_error() = 1/2`, so half epsilon is the maximum potential error.

```
double round_err = std::numeric_limits<double>::epsilon() // 2.2204460492503131e-016
                  * std::numeric_limits<double>::round_error(); // 1/2
std::cout << round_err << std::endl; // 1.1102230246251565e-016
```

There are, of course, many occasions when much bigger loss of precision occurs, for example, caused by [Loss of significance or cancellation error](#) or very many iterations.

## epsilon

Function `std::numeric_limits<T>::epsilon()` is meaningful only for non-integral types.

It returns the difference between 1.0 and the next value representable by the floating-point type T. So it is a one least-significant-bit change in this floating-point value.

For double (float\_64t) it is 2.2204460492503131e-016 showing all possibly significant 17 decimal digits.

```
std::cout.precision(std::numeric_limits<double>::max_digits10);
double d = 1.;
double eps = std::numeric_limits<double>::epsilon();
double dpeps = d+eps;
std::cout << std::showpoint // Ensure all trailing zeros are shown.
  << d << "\n"           // 1.0000000000000000
  << dpeps << std::endl; // 2.2204460492503131e-016
std::cout << dpeps - d    // 1.0000000000000002
  << std::endl;
```

We can explicitly increment by one bit using the function `boost::math::float_next()` and the result is the same as adding epsilon.

```
double one = 1.;
double nad = boost::math::float_next(one);
std::cout << nad << "\n" // 1.0000000000000002
  << nad - one // 2.2204460492503131e-016
  << std::endl;
```

Adding any smaller value, like half epsilon, will have no effect on this value.

```
std::cout.precision(std::numeric_limits<double>::max_digits10);
double d = 1.;
double eps = std::numeric_limits<double>::epsilon();
double dpeps = d + eps/2;

std::cout << std::showpoint // Ensure all trailing zeros are shown.
  << dpeps << "\n"           // 1.0000000000000000
  << eps/2 << std::endl; // 1.1102230246251565e-016
std::cout << dpeps - d    // 0.0000000000000000
  << std::endl;
```

So this cancellation error leaves the values equal, despite adding half epsilon.

To achieve greater portability over platform and floating-point type, Boost.Math and Boost.Multiprecision provide a package of functions that 'do something sensible' if the standard `numeric_limits` is not available. To use these `#include <boost/math/tools/precision.hpp>`.

A tolerance might be defined using this version of epsilon thus:

```
RealType tolerance = boost::math::tools::epsilon<RealType>() * 2;
```

## Tolerance for Floating-point Comparisons

epsilon is very useful to compute a tolerance when comparing floating-point values, a much more difficult task than is commonly imagined.

For more information you probably want (but still need) see [What Every Computer Scientist Should Know About Floating-Point Arithmetic](#)

The naive test comparing the absolute difference between two values and a tolerance does not give useful results if the values are too large or too small.

So Boost.Test uses an algorithm first devised by Knuth for reliably checking if floating-point values are close enough.

See Donald. E. Knuth. The art of computer programming (vol II). Copyright 1998 Addison-Wesley Longman, Inc., 0-201-89684-2. Addison-Wesley Professional; 3rd edition.



See also:

[Alberto Squassia, Comparing floats](#)

[Alberto Squassia, Comparing floats code](#)

[floating-point comparison](#).

For example, if we want a tolerance that might suit about 9 arithmetical operations, say  $\sqrt{9} = 3$ , we could define:

```
T tolerance = 3 * std::numeric_limits<T>::epsilon();
```

This is very widely used in Boost.Math testing with Boost.Test's macro `BOOST_CHECK_CLOSE_FRACTION`

```
T expected = 1.0;
T calculated = 1.0 + std::numeric_limits<T>::epsilon();

BOOST_CHECK_CLOSE_FRACTION(expected, calculated, tolerance);
```

used thus:

```
BOOST_CHECK_CLOSE_FRACTION(expected, calculated, tolerance);
```

(There is also a version using tolerance as a percentage rather than a fraction).

```
using boost::multiprecision::number;
using boost::multiprecision::cpp_dec_float;
using boost::multiprecision::et_off;

typedef number<cpp_dec_float<50>, et_off > cpp_dec_float_50; // 50 decimal digits.
```



## Note

that Boost.Test does not yet allow floating-point comparisons with expression templates on, so the default expression template parameter has been replaced by `et_off`.

```
cpp_dec_float_50 tolerance = 3 * std::numeric_limits<cpp_dec_float_50>::epsilon();
cpp_dec_float_50 expected = boost::math::constants::two_pi<cpp_dec_float_50>();
cpp_dec_float_50 calculated = 2 * boost::math::constants::pi<cpp_dec_float_50>();

BOOST_CHECK_CLOSE_FRACTION(expected, calculated, tolerance);
```

## Infinity - positive and negative

For floating-point types only, for which `std::numeric_limits<T>::has_infinity == true`, function `std::numeric_limits<T>::infinity()` provides an implementation-defined representation for  $\infty$ .

The 'representation' is a particular bit pattern reserved for infinity. For IEEE754 system (for which `std::numeric_limits<T>::is_iec559 == true`) [positive and negative infinity](#) are assigned bit patterns for all defined floating-point types.

Confusingly, the string resulting from outputting this representation, is also implementation-defined. And the string that can be input to generate the representation is also implementation-defined.

For example, the output is `1.#INF` on Microsoft systems, but `inf` on most \*nix platforms.

This implementation-defined-ness has hampered use of infinity (and NaNs) but Boost.Math and Boost.Multiprecision work hard to provide a sensible representation for **all** floating-point types, not just the built-in types, which with the use of suitable facets to define the input and output strings, makes it possible to use these useful features portably and including Boost.Serialization.

## Not-A-Number NaN

### Quiet\_NaN

For floating-point types only, for which `std::numeric_limits<T>::has_quiet_NaN == true`, function `std::numeric_limits<T>::quiet_NaN()` provides an implementation-defined representation for NaN.

NaNs are values to indicate that the result of an assignment or computation is meaningless. A typical example is `0/0` but there are many others.

NaNs may also be used, to represent missing values: for example, these could, by convention, be ignored in calculations of statistics like means.

Many of the problems with a representation for **Not-A-Number** has hampered portable use, similar to those with infinity.

NaN can be used with binary multiprecision types like `cpp_bin_float_quad`:

```
using boost::multiprecision::cpp_bin_float_quad;

if (std::numeric_limits<cpp_bin_float_quad>::has_quiet_NaN == true)
{
    cpp_bin_float_quad tolerance = 3 * std::numeric_limits<cpp_bin_float_quad>::epsilon();

    cpp_bin_float_quad NaN = std::numeric_limits<cpp_bin_float_quad>::quiet_NaN();
    std::cout << "cpp_bin_float_quad NaN is " << NaN << std::endl; // cpp_bin_float_quad NaN ↴
    is nan

    cpp_bin_float_quad expected = NaN;
    cpp_bin_float_quad calculated = 2 * NaN;
    // Comparisons of NaN's always fail:
    bool b = expected == calculated;
    std::cout << b << std::endl;
    BOOST_CHECK_NE(expected, expected);
    BOOST_CHECK_NE(expected, calculated);
}
else
{
    std::cout << "Type " << typeid(cpp_bin_float_quad).name() << " does not have NaNs!" << std::endl;
}
```

But using Boost.Math and suitable facets can permit portable use of both NaNs and positive and negative infinity.

See [boost:/libs/math/example/nonfinite\\_facet\\_sstream.cpp](http://boost.org/libs/math/example/nonfinite_facet_sstream.cpp) and we also need

```
#include <boost/math/special_functions/nonfinite_num_facets.hpp>
```

Then we can equally well use a multiprecision type `cpp_bin_float_quad`:

```
using boost::multiprecision::cpp_bin_float_quad;

typedef cpp_bin_float_quad T;

using boost::math::nonfinite_num_put;
using boost::math::nonfinite_num_get;
{
    std::locale old_locale;
    std::locale tmp_locale(old_locale, new nonfinite_num_put<char>);
    std::locale new_locale(tmp_locale, new nonfinite_num_get<char>);
    std::stringstream ss;
    ss.imbue(new_locale);
    T inf = std::numeric_limits<T>::infinity();
    ss << inf; // Write out.
    assert(ss.str() == "inf");
    T r;
    ss >> r; // Read back in.
    assert(inf == r); // Confirms that the floating-point values really are identical.
    std::cout << "infinity output was " << ss.str() << std::endl;
    std::cout << "infinity input was " << r << std::endl;
}
```

```
infinity output was inf
infinity input was inf
```

Similarly we can do the same with NaN (except that we cannot use assert)

```
{
    std::locale old_locale;
    std::locale tmp_locale(old_locale, new nonfinite_num_put<char>);
    std::locale new_locale(tmp_locale, new nonfinite_num_get<char>);
    std::stringstream ss;
    ss.imbue(new_locale);
    T n;
    T NaN = std::numeric_limits<T>::quiet_NaN();
    ss << NaN; // Write out.
    assert(ss.str() == "nan");
    std::cout << "NaN output was " << ss.str() << std::endl;
    ss >> n; // Read back in.
    std::cout << "NaN input was " << n << std::endl;
}
```

NaN output was nan NaN input was nan

## Signaling NaN

For floating-point types only, for which `std::numeric_limits<T>::has_signaling_NaN == true`, function `std::numeric_limits<T>::signaling_NaN()` provides an implementation-defined representation for NaN that causes a hardware trap. It should be noted however, that at least one implementation of this function causes a hardware trap to be triggered simply by calling `std::numeric_limits<T>::signaling_NaN()`, and not only by using the value returned.

## Numeric limits for 32-bit platform

These tables were generated using the following program and options:

```
Program:
  numeric_limits_qbk.cpp
Mon Nov  4 18:09:06 2013
BuildInfo:
  Platform Win32
  Compiler Microsoft Visual C++ version 10.0
  MSVC version 160040219.
  STL Dinkumware standard library version 520
  Boost version 1.55.0
```

**Table 4. Integer types constants** (`std::numeric_limits<T>::is_integer == true && is_exact == true`)

type	signed	bound	modulo	round	radix	digits	digits10
bool	unsigned	bound	no	to zero	2	1	0
char	signed	bound	modulo	to zero	2	7	2
unsigned char	unsigned	bound	modulo	to zero	2	8	2
char16_t	unsigned	bound	modulo	to zero	2	16	4
char32_t	unsigned	bound	modulo	to zero	2	32	9
short	signed	bound	modulo	to zero	2	15	4
unsigned short	unsigned	bound	modulo	to zero	2	16	4
int	signed	bound	modulo	to zero	2	31	9
unsigned	unsigned	bound	modulo	to zero	2	32	9
long	signed	bound	modulo	to zero	2	31	9
unsigned long	unsigned	bound	modulo	to zero	2	32	9
long long	signed	bound	modulo	to zero	2	63	18
unsigned long long	unsigned	bound	modulo	to zero	2	64	19
int32_t	signed	bound	modulo	to zero	2	31	9
uint32_t	unsigned	bound	modulo	to zero	2	32	9
int64_t	signed	bound	modulo	to zero	2	63	18
uint64_t	unsigned	bound	modulo	to zero	2	64	19
int128_t	signed	bound	modulo	to zero	2	128	38
uint128_t	unsigned	bound	modulo	to zero	2	128	38
int256_t	signed	bound	modulo	to zero	2	256	77
uint256_t	unsigned	bound	modulo	to zero	2	256	77
cpp_int	signed	unbounded	no	to zero	2	2147483647	646392383

**Table 5. Integer types functions** (`std::numeric_limits<T>::is_integer == true` && `std::numeric_limits<T>::min() == std::numeric_limits<T>::lowest()`)

function	max	min
bool	1	0
char	127	-128
unsigned char	255	0
char16_t	65535	0
char32_t	4294967295	0
short	32767	-32768
unsigned short	65535	0
int	2147483647	-2147483648
unsigned int	4294967295	0
long	2147483647	-2147483648
unsigned long	4294967295	0
long long	9223372036854775807	-9223372036854775808
unsigned long long	18446744073709551615	0
int32_t	2147483647	-2147483648
int64_t	9223372036854775807	-9223372036854775808
int128_t	340282366920938463463374607431768211455	-340282366920938463463374607431768211455

**Table 6. Floating-point types constants** (`std::numeric_limits<T>::is_integer==false && is_signed==true && is_modulo==false && is_exact==false && is_bound==true`)

type	round	radix	digits	digits10	max_digits10	min_exp	min_exp10	max_exp	max_exp10	tiny	trap
float	to nearest	2	24	6	8	-125	-37	128	38	tiny	traps
double	to nearest	2	53	15	17	-1021	-307	1024	308	tiny	traps
long double	to nearest	2	53	15	17	-1021	-307	1024	308	tiny	traps
cpp_dec_50	indeterminate	10	50	50	80	-225300	-67108864	225300	67108864	no	no
cpp_bin_128	to nearest	2	377	113	115	-2478288	-646456766	2478288	646456766	no	traps

**Table 7. Floating-point types functions** (`std::numeric_limits<T>::is_integer == false`)

function	float	double	long double	cpp_dec_50	cpp_bin_128
max	3.40282e+038	1.79769e+308	1.79769e+308	1e+67108865	1.85906e+646456766
min	1.17549e-038	2.22507e-308	2.22507e-308	1e-67108864	5.37906e-646456767
epsilon	1.19209e-007	2.22045e-016	2.22045e-016	1e-49	6.49713e-114
round_error	0.5	0.5	0.5	0.5	0.5
infinity	1.#INF	1.#INF	1.#INF	inf	inf
quiet_NaN	1.#QNAN	1.#QNAN	1.#QNAN	nan	nan
signaling_NaN	1.#QNAN	1.#QNAN	1.#QNAN	0	0
denorm_min	1.4013e-045	4.94066e-324	4.94066e-324	0	0

## How to Determine the Kind of a Number From `std::numeric_limits`

Based on the information above, one can see that different kinds of numbers can be differentiated based on the information stored in `std::numeric_limits`. This is in addition to the traits class `number_category` provided by this library.

## Integer Types

For an integer type T, all of the following conditions hold:

```
std::numeric_limits<T>::is_specialized == true
std::numeric_limits<T>::is_integer == true
std::numeric_limits<T>::is_exact == true
std::numeric_limits<T>::min_exponent == 0
std::numeric_limits<T>::max_exponent == 0
std::numeric_limits<T>::min_exponent10 == 0
std::numeric_limits<T>::max_exponent10 == 0
```

In addition the type is *signed* if:

```
std::numeric_limits<T>::is_signed == true
```

If the type is arbitrary precision then:

```
std::numeric_limits<T>::is_bounded == false
```

Otherwise the type is bounded, and returns a non zero value from:

```
std::numeric_limits<T>::max()
```

and has:

```
std::numeric_limits<T>::is_modulo == true
```

if the type implements modulo arithmetic on overflow.

## Rational Types

Rational types are just like integers except that:

```
std::numeric_limits<T>::is_integer == false
```

## Fixed Precision Types

There appears to be no way to tell these apart from rational types, unless they set:

```
std::numeric_limits<T>::is_exact == false
```

This is because these types are in essence a rational type with a fixed denominator.

## Floating Point Types

For a floating point type T, all of the following conditions hold:

```
std::numeric_limits<T>::is_specialized == true
std::numeric_limits<T>::is_integer == false
std::numeric_limits<T>::is_exact == false
std::numeric_limits<T>::min_exponent != 0
std::numeric_limits<T>::max_exponent != 0
std::numeric_limits<T>::min_exponent10 != 0
std::numeric_limits<T>::max_exponent10 != 0
```

In addition the type is *signed* if:



```
std::numeric_limits<T>::is_signed == true
```

And the type may be decimal or binary depending on the value of:

```
std::numeric_limits<T>::radix
```

In general, there are no arbitrary precision floating point types, and so:

```
std::numeric_limits<T>::is_bounded == false
```

## Exact Floating Point Types

Exact floating point types are a [field](#) composed of an arbitrary precision integer scaled by an exponent. Such types have no division operator and are the same as floating point types except:

```
std::numeric_limits<T>::is_exact == true
```

## Complex Numbers

For historical reasons, complex numbers do not specialize `std::numeric_limits`, instead you must inspect `std::numeric_limits<T::value_type>`.

# Input Output

## Loopback testing

*Loopback* or *round-tripping* refers to writing out a value as a decimal digit string using `std::ostream`, usually to a `std::stringstream`, and then reading the string back in to another value, and confirming that the two values are identical. A trivial example using `float` is:

```
float write; // Value to round-trip.
std::stringstream ss; // Read and write std::stringstream.
ss.precision(std::numeric_limits<T>::max_digits10); // Ensure all potentially significant bits ↴
are output.
ss.flags(std::ios_base::fmtflags(std::ios_base::scientific)); // Use scientific format.
ss << write; // Output to string.
float read; // Expected.
ss >> read; // Read decimal digits string from stringstream.
BOOST_CHECK_EQUAL(write, read); // Should be the same.
```

and this can be run in a loop for all possible values of a 32-bit float. For other floating-point types `T`, including built-in `double`, it takes far too long to test all values, so a reasonable test strategy is to use a large number of random values.

```

T write;
std::stringstream ss;
ss.precision(std::numeric_limits<T>::max_digits10); // Ensure all potentially significant bits
are output.
ss.flags(f); // Changed from default iostream format flags if desired.
ss << write; // Output to stringstream.

T read;
ss >> read; // Get read using operator>> from stringstream.
BOOST_CHECK_EQUAL(read, write);

read = static_cast<T>(ss.str()); // Get read by converting from decimal digits string represent-
ation of write.
BOOST_CHECK_EQUAL(read, write);

read = static_cast<T>(write.str(0, f)); // Get read using format specified when written.
BOOST_CHECK_EQUAL(read, write);

```

The test at [test\\_cpp\\_bin\\_float\\_io.cpp](#) allows any floating-point type to be *round\_tripped* using a wide range of fairly random values. It also includes tests compared a collection of [stringdata](#) test cases in a file.

## Comparing with output using Built-in types

One can make some comparisons with the output of

```
<number<cpp_bin_float<53, digit_count_2> >
```

which has the same number of significant bits (53) as 64-bit double precision floating-point.

However, although most outputs are identical, there are differences on some platforms caused by the implementation-dependent behaviours allowed by the C99 specification [C99 ISO/IEC 9899:TC2](#), incorporated by C++.

*"For e, E, f, F, g, and G conversions, if the number of significant decimal digits is at most DECIMAL\_DIG, then the result should be correctly rounded. If the number of significant decimal digits is more than DECIMAL\_DIG but the source value is exactly representable with DECIMAL\_DIG digits, then the result should be an exact representation with trailing zeros. Otherwise, the source value is bounded by two adjacent decimal strings  $L < U$ , both having DECIMAL\_DIG significant digits; the value of the resultant decimal string  $D$  should satisfy  $L \leq D \leq U$ , with the extra stipulation that the error should have a correct sign for the current rounding direction."*

So not only is correct rounding for the full number of digits not required, but even if the **optional** recommended practice is followed, then the value of these last few digits is unspecified as long as the value is within certain bounds.



### Note

Do not expect the output from different platforms to be **identical**, but `cpp_dec_float`, `cpp_bin_float` (and other backends) outputs should be correctly rounded to the number of digits requested by the set precision and format.

## Macro BOOST\_MP\_MIN\_EXPONENT\_DIGITS

[C99 Standard](#) for format specifiers, 7.19.6 Formatted input/output functions requires:

"The exponent always contains at least two digits, and only as many more digits as necessary to represent the exponent."

So to conform to the C99 standard (incorporated by C++)

```
#define BOOST_MP_MIN_EXPONENT_DIGITS 2
```

Confusingly, Microsoft (and MinGW) do not conform to this standard and provide **at least three digits**, for example `1e+001`. So if you want the output to match that from built-in floating-point types on compilers that use Microsofts runtime then use:

```
#define BOOST_MP_MIN_EXPONENT_DIGITS 3
```

Also useful to get the minimum exponent field width is

```
#define BOOST_MP_MIN_EXPONENT_DIGITS 1
```

producing a compact output like `2e+4`, useful when conserving space is important.

Larger values are also supported, for example, value 4 for `2e+0004` which may be useful to ensure that columns line up.

# Reference

## number

### Synopsis

```

namespace boost{ namespace multiprecision{

enum expression_template_option { et_on = 1, et_off = 0 };

template <class Backend> struct expression_template_default
{ static const expression_template_option value = et_on; };

template <class Backend, expression_template_option ExpressionTemplates = expression_template_de
fault<Backend>::value>
class number
{
    number();
    number(see-below);
    number& operator=(see-below);
    number& assign(see-below);

    // Member operators
    number& operator+=(const see-below&);
    number& operator-=(const see-below&);
    number& operator*=(const see-below&);
    number& operator/=(const see-below&);
    number& operator++();
    number& operator--();
    number operator++(int);
    number operator--(int);

    number& operator%=(const see-below&);
    number& operator&=(const see-below&);
    number& operator|=(const see-below&);
    number& operator^=(const see-below&);
    number& operator<=<=(const integer-type&);
    number& operator>=>=(const integer-type&);

    // Use in Boolean context:
    operator convertible-to-bool-type()const;
    // swap:
    void swap(number& other);
    // Sign:
    bool is_zero()const;
    int sign()const;
    // string conversion:
    std::string str()const;
    // Generic conversion mechanism
    template <class T>
    T convert_to()const;
    template <class T>
    explicit operator T ()const;
    // precision control:
    static unsigned default_precision();
    static void default_precision(unsigned digits10);
    unsigned precision()const;
    void precision(unsigned digits10);
    // Comparison:
    int compare(const number<Backend>& o)const;
    template <class V>

```

```

typename enable_if<is_convertible<V, number<Backend, ExpressionTemplates> >, int>::type
    compare(const V& o) const;
// Access to the underlying implementation:
Backend& backend();
const Backend& backend() const;
};

// Non member operators:
unmentionable-expression-template-type operator+(const see-below&);
unmentionable-expression-template-type operator-(const see-below&);
unmentionable-expression-template-type operator+(const see-below&, const see-below&);
unmentionable-expression-template-type operator-(const see-below&, const see-below&);
unmentionable-expression-template-type operator*(const see-below&, const see-below&);
unmentionable-expression-template-type operator/(const see-below&, const see-below&);
// Integer only operations:
unmentionable-expression-template-type operator%(const see-below&, const see-below&);
unmentionable-expression-template-type operator&(const see-below&, const see-below&);
unmentionable-expression-template-type operator|(const see-below&, const see-below&);
unmentionable-expression-template-type operator^(const see-below&, const see-below&);
unmentionable-expression-template-type operator<<(const see-below&, const integer-type&);
unmentionable-expression-template-type operator>>(const see-below&, const integer-type&);
// Comparison operators:
bool operator==(const see-below&, const see-below&);
bool operator!=(const see-below&, const see-below&);
bool operator<(const see-below&, const see-below&);
bool operator>(const see-below&, const see-below&);
bool operator<=(const see-below&, const see-below&);
bool operator>=(const see-below&, const see-below&);

// Swap:
template <class Backend, expression_template_option ExpressionTemplates>
void swap(number<Backend, ExpressionTemplates>& a, number<Backend, ExpressionTemplates>& b);

// iostream support:
template <class Backend, expression_template_option ExpressionTemplates>
std::ostream& operator << (std::ostream& os, const number<Backend, ExpressionTemplates>& r);
std::ostream& operator << (std::ostream& os, const unmentionable-expression-template-type& r);
template <class Backend, expression_template_option ExpressionTemplates>
std::istream& operator >> (std::istream& is, number<Backend, ExpressionTemplates>& r);

// Arithmetic with a higher precision result:
template <class ResultType, class Source1 class Source2>
ResultType& add(ResultType& result, const Source1& a, const Source2& b);
template <class ResultType, class Source1 class Source2>
ResultType& subtract(ResultType& result, const Source1& a, const Source2& b);
template <class ResultType, class Source1 class Source2>
ResultType& multiply(ResultType& result, const Source1& a, const Source2& b);

// Non-member function standard library support:
unmentionable-expression-template-type abs (const number-or-expression-template-type&);
unmentionable-expression-template-type fabs (const number-or-expression-template-type&);
unmentionable-expression-template-type sqrt (const number-or-expression-template-type&);
unmentionable-expression-template-type floor (const number-or-expression-template-type&);
unmentionable-expression-template-type ceil (const number-or-expression-template-type&);
unmentionable-expression-template-type trunc (const number-or-expression-template-type&);
unmentionable-expression-template-type itrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type ltrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type lltrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type round (const number-or-expression-template-type&);
unmentionable-expression-template-type iround (const number-or-expression-template-type&);
unmentionable-expression-template-type lround (const number-or-expression-template-type&);
unmentionable-expression-template-type llround (const number-or-expression-template-type&);
unmentionable-expression-template-type exp (const number-or-expression-template-type&);

```

```

unmentionable-expression-template-type    log      (const number-or-expression-template-type&);
unmentionable-expression-template-type    log10    (const number-or-expression-template-type&);
unmentionable-expression-template-type    cos      (const number-or-expression-template-type&);
unmentionable-expression-template-type    sin      (const number-or-expression-template-type&);
unmentionable-expression-template-type    tan      (const number-or-expression-template-type&);
unmentionable-expression-template-type    acos     (const number-or-expression-template-type&);
unmentionable-expression-template-type    asin     (const number-or-expression-template-type&);
unmentionable-expression-template-type    atan     (const number-or-expression-template-type&);
unmentionable-expression-template-type    cosh     (const number-or-expression-template-type&);
unmentionable-expression-template-type    sinh     (const number-or-expression-template-type&);
unmentionable-expression-template-type    tanh     (const number-or-expression-template-type&);

unmentionable-expression-template-type    ldexp    (const number-or-expression-template-type&, int);
unmentionable-expression-template-type    frexp    (const number-or-expression-template-type&, int*);
unmentionable-expression-template-type    pow      (const number-or-expression-template-
type&, const number-or-expression-template-type&);
unmentionable-expression-template-type    fmod     (const number-or-expression-template-
type&, const number-or-expression-template-type&);
unmentionable-expression-template-type    atan2    (const number-or-expression-template-
type&, const number-or-expression-template-type&);

// Traits support:
template <class T>
struct component_type;
template <class T>
struct number_category;
template <class T>
struct is_number;
template <class T>
struct is_number_expression;

// Integer specific functions:
unmentionable-expression-template-type    gcd(const number-or-expression-template-
type&, const number-or-expression-template-type&);
unmentionable-expression-template-type    lcm(const number-or-expression-template-
type&, const number-or-expression-template-type&);
unmentionable-expression-template-type    pow(const number-or-expression-template-type&, unsigned);
unmentionable-expression-template-type    powm(const number-or-expression-template-
type& b, const number-or-expression-template-type& p, const number-or-expression-template-type& m);
unmentionable-expression-template-type    sqrt(const number-or-expression-template-type&);
template <class Backend, expression_template_option ExpressionTemplates>
number<Backend, ExpressionTemplates>      sqrt(const number-or-expression-template-type&, num-
ber<Backend, ExpressionTemplates>&);
template <class Backend, expression_template_option ExpressionTemplates>
void divide_qr(const number-or-expression-template-type& x, const number-or-expression-template-
type& y,
               number<Backend, ExpressionTemplates>& q, number<Backend, ExpressionTemplates>& r);
template <class Integer>
Integer integer_modulus(const number-or-expression-template-type& x, Integer val);
unsigned lsb(const number-or-expression-template-type& x);
unsigned msb(const number-or-expression-template-type& x);
template <class Backend, class ExpressionTemplates>
bool bit_test(const number<Backend, ExpressionTemplates>& val, unsigned index);
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_set(number<Backend, ExpressionTemplates>& val, un-
signed index);
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_unset(number<Backend, ExpressionTemplates>& val, un-
signed index);
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_flip(number<Backend, ExpressionTemplates>& val, un-
signed index);
template <class Engine>

```

```

bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials, Engine& gen);
bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials);

// Rational number support:
typename component_type<number-or-expression-template-type>::type numerator (const number-or-
expression-template-type&);
typename component_type<number-or-expression-template-type>::type denominator(const number-or-
expression-template-type&);

}} // namespaces

namespace boost{ namespace math{

// Boost.Math interoperability functions:
int fpclassify (const number-or-expression-
template-type&, int);
bool isfinite (const number-or-expression-
template-type&, int);
bool isnan (const number-or-expression-
template-type&, int);
bool isinf (const number-or-expression-
template-type&, int);
bool isnormal (const number-or-expression-
template-type&, int);

}} // namespaces

// numeric_limits support:
namespace std{

template <class Backend, expression_template_option ExpressionTemplates>
struct numeric_limits<boost::multiprecision<Backend, ExpressionTemplates> >
{
    /* Usual members here */
};

}

```

## Description

```
enum expression_template_option { et_on = 1, et_off = 0 };
```

This enumerated type is used to specify whether expression templates are turned on (et\_on) or turned off (et\_off).

```
template <class Backend> struct expression_template_default
{ static const expression_template_option value = et_on; };
```

This traits class specifies the default expression template option to be used with a particular Backend type. It defaults to et\_on.

```
template <class Backend, expression_template_option ExpressionTemplates = expression_template_de-
fault<Backend>::value>
class number;
```

Class number has two template arguments:

Backend	The actual arithmetic back-end that does all the work.
ExpressionTemplates	A Boolean value: when et_on, then expression templates are enabled, otherwise when set to et_off they are disabled. The default for this parameter is computed via the traits class expression_tem-

`plate_default` whose member value defaults to `et_on` unless the traits class is specialized for a particular backend.

```
number();  
number(see-below);  
number& operator=(see-below);  
number& assign(see-below);
```

Type `number` is default constructible, and both copy constructible and assignable from:

- Itself.
- An expression template which is the result of one of the arithmetic operators.
- Any builtin arithmetic type, as long as the result would not be lossy (for example float to integer conversion).
- Any type that the Backend is implicitly constructible or assignable from.
- An rvalue reference to another `number`. Move-semantics are used for construction if the backend also supports rvalue reference construction. In the case of assignment, move semantics are always supported when the argument is an rvalue reference irrespective of the backend.
- Any type in the same family, as long as no loss of precision is involved. For example from `int128_t` to `int256_t`, or `cpp_dec_float_50` to `cpp_dec_float_100`.

Type `number` is explicitly constructible from:

- Any type mentioned above.
- A `std::string` or any type which is convertible to `const char*`.
- Any arithmetic type (including those that would result in lossy conversions).
- Any type in the same family, including those that result in loss of precision.
- Any type that the Backend is explicitly constructible from.
- Any pair of types for which a generic interconversion exists: that is from integer to integer, integer to rational, integer to float, rational to rational, rational to float, or float to float.

The `assign` member function is available for any type for which an explicit converting constructor exists. It is intended to be used where a temporary generated from an explicit assignment would be expensive, for example:

```
mpfr_float_50    f50;  
mpfr_float_100   f100;  
  
f50 = static_cast<mpfr_float_50>(f100); // explicit cast create a temporary  
f50.assign(f100);                       // explicit call to assign create no temporary
```

In addition, if the type has multiple components (for example rational or complex number types), then there is a two argument constructor:

```
number(arg1, arg2);
```

Where the two args must either be arithmetic types, or types that are convertible to the two components of `this`.



```
number& operator+=(const see-below&);
number& operator-=(const see-below&);
number& operator*=(const see-below&);
number& operator/=(const see-below&);
number& operator++();
number& operator--();
number operator++(int);
number operator--(int);
// Integer only operations:
number& operator%=(const see-below&);
number& operator&=(const see-below&);
number& operator|=(const see-below&);
number& operator^=(const see-below&);
number& operator<=<=(const integer-type&);
number& operator>=>=(const integer-type&);
```

These operators all take their usual arithmetic meanings.

The arguments to these operators is either:

- Another `number<Backend, ExpressionTemplates>`.
- An expression template derived from `number<Backend>`.
- Any type implicitly convertible to `number<Backend, ExpressionTemplates>`, including some other instance of class `number`.

For the left and right shift operations, the argument must be a builtin integer type with a positive value (negative values result in a `std::runtime_error` being thrown).

```
operator convertible-to-bool-type()const;
```

Returns an *unmentionable-type* that is usable in Boolean contexts (this allows `number` to be used in any Boolean context - if statements, conditional statements, or as an argument to a logical operator - without type `number` being convertible to type `bool`).

This operator also enables the use of `number` with any of the following operators: `!`, `|`, `&&` and `? :`.

```
void swap(number& other);
```

Swaps `*this` with `other`.

```
bool is_zero()const;
```

Returns `true` if `*this` is zero, otherwise `false`.

```
int sign()const;
```

Returns a value less than zero if `*this` is negative, a value greater than zero if `*this` is positive, and zero if `*this` is zero.

```
std::string str(unsigned precision, bool scientific = true)const;
```

Returns the number formatted as a string, with at least *precision* digits, and in scientific format if *scientific* is `true`.

```
template <class T>
T convert_to()const;

template <class T>
explicit operator T ()const;
```

Provides a generic conversion mechanism to convert `*this` to type `T`. Type `T` may be any arithmetic type. Optionally other types may also be supported by specific Backend types.

```
static unsigned default_precision();
static void default_precision(unsigned digits10);
unsigned precision()const;
void precision(unsigned digits10);
```

These functions are only available if the Backend template parameter supports runtime changes to precision. They get and set the default precision and the precision of `*this` respectively.

```
int compare(const number<Backend, ExpressionTemplates>& o)const;
template <class V>
typename enable_if<is_convertible<V, number<Backend, ExpressionTemplates> >, int>::type
compare(const V& other)const;
```

Returns:

- A value less than 0 for `*this < other`
- A value greater than 0 for `*this > other`
- Zero for `*this == other`

```
Backend& backend();
const Backend& backend()const;
```

Returns the underlying back-end instance used by `*this`.

## Non-member operators

```
// Non member operators:
unmentionable-expression-template-type operator+(const see-below&);
unmentionable-expression-template-type operator-(const see-below&);
unmentionable-expression-template-type operator+(const see-below&, const see-below&);
unmentionable-expression-template-type operator-(const see-below&, const see-below&);
unmentionable-expression-template-type operator*(const see-below&, const see-below&);
unmentionable-expression-template-type operator/(const see-below&, const see-below&);
// Integer only operations:
unmentionable-expression-template-type operator%(const see-below&, const see-below&);
unmentionable-expression-template-type operator&(const see-below&, const see-below&);
unmentionable-expression-template-type operator|(const see-below&, const see-below&);
unmentionable-expression-template-type operator^(const see-below&, const see-below&);
unmentionable-expression-template-type operator<<(const see-below&, const integer-type&);
unmentionable-expression-template-type operator>>(const see-below&, const integer-type&);
// Comparison operators:
bool operator==(const see-below&, const see-below&);
bool operator!=(const see-below&, const see-below&);
bool operator<(const see-below&, const see-below&);
bool operator>(const see-below&, const see-below&);
bool operator<=(const see-below&, const see-below&);
bool operator>=(const see-below&, const see-below&);
```

These operators all take their usual arithmetic meanings.

The arguments to these functions must contain at least one of the following:

- A number.
- An expression template type derived from number.
- Any type for which number has an implicit constructor - for example a builtin arithmetic type.

The return type of these operators is either:

- An *unmentionable-type* expression template type when `ExpressionTemplates` is true.
- Type `number<Backend, et_off>` when `ExpressionTemplates` is false.
- Type `bool` if the operator is a comparison operator.

Finally note that the second argument to the left and right shift operations must be a builtin integer type, and that the argument must be positive (negative arguments result in a `std::runtime_error` being thrown).

## swap

```
template <class Backend, ExpressionTemplates>
void swap(number<Backend, ExpressionTemplates>& a, number<Backend, ExpressionTemplates>& b);
```

Swaps a and b.

## ostream Support

```
template <class Backend, expression_template_option ExpressionTemplates>
std::ostream& operator << (std::ostream& os, const number<Backend, ExpressionTemplates>& r);
template <class Unspecified...>
std::ostream& operator << (std::ostream& os, const unmentionable-expression-template& r);
template <class Backend, expression_template_option ExpressionTemplates>
inline std::istream& operator >> (std::istream& is, number<Backend, ExpressionTemplates>& r)
```

These operators provided formatted input-output operations on number types, and expression templates derived from them.

It's down to the back-end type to actually implement string conversion. However, the back-ends provided with this library support all of the iostream formatting flags, field width and precision settings.

## Arithmetic with a higher precision result

```
template <class ResultType, class Source1 class Source2>
ResultType& add(ResultType& result, const Source1& a, const Source2& b);

template <class ResultType, class Source1 class Source2>
ResultType& subtract(ResultType& result, const Source1& a, const Source2& b);

template <class ResultType, class Source1 class Source2>
ResultType& multiply(ResultType& result, const Source1& a, const Source2& b);
```

These functions apply the named operator to the arguments *a* and *b* and store the result in *result*, returning *result*. In all cases they behave "as if" arguments *a* and *b* were first promoted to type *ResultType* before applying the operator, though particular backends may well avoid that step by way of an optimization.

The type *ResultType* must be an instance of class *number*, and the types *Source1* and *Source2* may be either instances of class *number* or native integer types. The latter is an optimization that allows arithmetic to be performed on native integer types producing an extended precision result.

## Non-member standard library function support

```
unmentionable-expression-template-type abs (const number-or-expression-template-type&);
unmentionable-expression-template-type fabs (const number-or-expression-template-type&);
unmentionable-expression-template-type sqrt (const number-or-expression-template-type&);
unmentionable-expression-template-type floor (const number-or-expression-template-type&);
unmentionable-expression-template-type ceil (const number-or-expression-template-type&);
unmentionable-expression-template-type trunc (const number-or-expression-template-type&);
unmentionable-expression-template-type itrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type ltrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type lltrunc (const number-or-expression-template-type&);
unmentionable-expression-template-type round (const number-or-expression-template-type&);
unmentionable-expression-template-type iround (const number-or-expression-template-type&);
unmentionable-expression-template-type lround (const number-or-expression-template-type&);
unmentionable-expression-template-type llround (const number-or-expression-template-type&);
unmentionable-expression-template-type exp (const number-or-expression-template-type&);
unmentionable-expression-template-type log (const number-or-expression-template-type&);
unmentionable-expression-template-type log10 (const number-or-expression-template-type&);
unmentionable-expression-template-type cos (const number-or-expression-template-type&);
unmentionable-expression-template-type sin (const number-or-expression-template-type&);
unmentionable-expression-template-type tan (const number-or-expression-template-type&);
unmentionable-expression-template-type acos (const number-or-expression-template-type&);
unmentionable-expression-template-type asin (const number-or-expression-template-type&);
unmentionable-expression-template-type atan (const number-or-expression-template-type&);
unmentionable-expression-template-type cosh (const number-or-expression-template-type&);
unmentionable-expression-template-type sinh (const number-or-expression-template-type&);
unmentionable-expression-template-type tanh (const number-or-expression-template-type&);

unmentionable-expression-template-type ldexp (const number-or-expression-template-type&, int);
unmentionable-expression-template-type frexp (const number-or-expression-template-type&, int*);
unmentionable-expression-template-type pow (const number-or-expression-template-type&, const number-or-expression-template-type&);
unmentionable-expression-template-type fmod (const number-or-expression-template-type&, const number-or-expression-template-type&);
unmentionable-expression-template-type atan2 (const number-or-expression-template-type&, const number-or-expression-template-type&);
```

These functions all behave exactly as their standard library C++11 counterparts do: their argument is either an instance of `number` or an expression template derived from it; If the argument is of type `number<Backend, et_off>` then that is also the return type, otherwise the return type is an expression template.

These functions are normally implemented by the Backend type. However, default versions are provided for Backend types that don't have native support for these functions. Please note however, that this default support requires the precision of the type to be a compile time constant - this means for example that the [GMP MPF Backend](#) will not work with these functions when that type is used at variable precision.

Also note that with the exception of `abs` that these functions can only be used with floating-point Backend types (if any other types such as fixed precision or complex types are added to the library later, then these functions may be extended to support those number types).

The precision of these functions is generally determined by the backend implementation. For example the precision of these functions when used with `mpfr_float` is determined entirely by [MPFR](#). When these functions use our own implementations, the accuracy of the transcendental functions is generally a few epsilon. Note however, that the trigonometrical functions incur the usual accuracy loss when reducing arguments by large multiples of  $\pi$ . Also note that both `gmp_float` and `cpp_dec_float` have a number of guard digits beyond their stated precision, so the error rates listed for these are in some sense artificially low.

The following table shows the error rates we observe for these functions with various backend types, functions not listed here are exact (tested on Win32 with VC++10, MPFR-3.0.0, MPIR-2.1.1):

Function	mpfr_float_50	mpf_float_50	cpp_dec_float_50
sqrt	1eps	0eps	0eps
exp	1eps	0eps	0eps
log	1eps	0eps	0eps
log10	1eps	0eps	0eps
cos	700eps	0eps	0eps
sin	1eps	0eps	0eps
tan	0eps	0eps	0eps
acos	0eps	0eps	0eps
asin	0eps	0eps	0eps
atan	1eps	0eps	0eps
cosh	1045eps <sup>1</sup>	0eps	0eps
sinh	2eps	0eps	0eps
tanh	1eps	0eps	0eps
pow	0eps	4eps	3eps
atan2	1eps	0eps	0eps

<sup>1</sup> It's likely that the inherent error in the input values to our test cases are to blame here.

## Traits Class Support

```
template <class T>
struct component_type;
```

If this is a type with multiple components (for example rational or complex types), then this trait has a single member `type` that is the type of those components.

```
template <class T>
struct number_category;
```

A traits class that inherits from `mpl::int_<N>` where `N` is one of the enumerated values `number_kind_integer`, `number_kind_floating_point`, `number_kind_rational`, `number_kind_fixed_point`, or `number_kind_unknown`. This traits class is specialized for any type that has `std::numeric_limits` support as well as for classes in this library: which means it can be used for generic code that must work with built in arithmetic types as well as multiprecision ones.

```
template <class T>
struct is_number;
```

A traits class that inherits from `mpl::true_` if `T` is an instance of `number<>`, otherwise from `mpl::false_`.

```
template <class T>
struct is_number_expression;
```

A traits class that inherits from `mpl::true_` if `T` is an expression template type derived from `number<>`, otherwise from `mpl::false_`.

## Integer functions

In addition to functioning with types from this library, these functions are also overloaded for built in integer types if you include `<boost/multiprecision/integer.hpp>`. Further, when used with fixed precision types (whether built in integers or multiprecision ones), the functions will promote to a wider type internally when the algorithm requires it. Versions overloaded for built in integer types return that integer type rather than an expression template.

```
unmentionable-expression-template-type gcd(const number-or-expression-template-
type& a, const number-or-expression-template-type& b);
```

Returns the largest integer `x` that divides both `a` and `b`.

```
unmentionable-expression-template-type lcm(const number-or-expression-template-
type& a, const number-or-expression-template-type& b);
```

Returns the smallest integer `x` that is divisible by both `a` and `b`.

```
unmentionable-expression-template-type pow(const number-or-expression-template-type& b, un-
signed p);
```

Returns  $b^p$  as an expression template. Note that this function should be used with extreme care as the result can grow so large as to take "effectively forever" to compute, or else simply run the host machine out of memory. This is the one function in this category that is not overloaded for built in integer types, further, it's probably not a good idea to use it with fixed precision `cpp_int`'s either.

```
unmentionable-expression-template-type powm(const number-or-expression-template-
type& b, const number-or-expression-template-type& p, const number-or-expression-template-type& m);
```

Returns  $b^p \bmod m$  as an expression template. Fixed precision types are promoted internally to ensure accuracy.

```
unmentionable-expression-template-type    sqrt(const number-or-expression-template-type& a);
```

Returns the largest integer  $x$  such that  $x * x < a$ .

```
template <class Backend, expression_template_option ExpressionTemplates>
number<Backend, ExpressionTemplates>    sqrt(const number-or-expression-template-
type& a, number<Backend, ExpressionTemplates>& r);
```

Returns the largest integer  $x$  such that  $x * x < a$ , and sets the remainder  $r$  such that  $r = a - x * x$ .

```
template <class Backend, expression_template_option ExpressionTemplates>
void divide_qr(const number-or-expression-template-type& x, const number-or-expression-template-
type& y,
              number<Backend, ExpressionTemplates>& q, number<Backend, ExpressionTemplates>& r);
```

Divides  $x$  by  $y$  and returns both the quotient and remainder. After the call  $q = x / y$  and  $r = x \% y$ .

```
template <class Integer>
Integer integer_modulus(const number-or-expression-template-type& x, Integer val);
```

Returns the absolute value of  $x \% val$ .

```
unsigned lsb(const number-or-expression-template-type& x);
```

Returns the (zero-based) index of the least significant bit that is set to 1.

Throws a `std::range_error` if the argument is  $\leq 0$ .

```
unsigned msb(const number-or-expression-template-type& x);
```

Returns the (zero-based) index of the most significant bit.

Throws a `std::range_error` if the argument is  $\leq 0$ .

```
template <class Backend, class ExpressionTemplates>
bool bit_test(const number<Backend, ExpressionTemplates>& val, unsigned index);
```

Returns true if the bit at *index* in *val* is set.

```
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_set(number<Backend, ExpressionTemplates>& val, un-
signed index);
```

Sets the bit at *index* in *val*, and returns *val*.

```
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_unset(number<Backend, ExpressionTemplates>& val, un-
signed index);
```

Unsets the bit at *index* in *val*, and returns *val*.

```
template <class Backend, class ExpressionTemplates>
number<Backend, ExpressionTemplates>& bit_flip(number<Backend, ExpressionTemplates>& val, un-
signed index);
```

Flips the bit at *index* in *val*, and returns *val*.

```
template <class Engine>
bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials, Engine& gen);
bool miller_rabin_test(const number-or-expression-template-type& n, unsigned trials);
```

Tests to see if the number *n* is probably prime - the test excludes the vast majority of composite numbers by excluding small prime factors and performing a single Fermat test. Then performs *trials* Miller-Rabin tests. Returns *false* if *n* is definitely composite, or *true* if *n* is probably prime with the probability of it being composite less than  $0.25^{\text{trials}}$ . Fixed precision types are promoted internally to ensure accuracy.

## Rational Number Functions

```
typename component_type<number-or-expression-template-type>::type numerator (const number-or-
expression-template-type&);
typename component_type<number-or-expression-template-type>::type denominator(const number-or-
expression-template-type&);
```

These functions return the numerator and denominator of a rational number respectively.

## Boost.Math Interoperability Support

```
namespace boost{ namespace math{

int fpclassify      (const number-or-expression-template-type&, int);
bool isfinite       (const number-or-expression-template-type&, int);
bool isnan          (const number-or-expression-template-type&, int);
bool isinf          (const number-or-expression-template-type&, int);
bool isnormal       (const number-or-expression-template-type&, int);

}} // namespaces
```

These floating-point classification functions behave exactly as their Boost.Math equivalents.

Other Boost.Math functions and templates may also be specialized or overloaded to ensure interoperability.

## std::numeric\_limits support

```
namespace std{

template <class Backend, ExpressionTemplates>
struct numeric_limits<boost::multiprecision<Backend, ExpressionTemplates> >
{
    /* Usual members here */
};

}
```

Class template `std::numeric_limits` is specialized for all instantiations of number whose precision is known at compile time, plus those types whose precision is unlimited (though it is much less useful in those cases). It is not specialized for types whose precision can vary at compile time (such as `mpf_float`).



## cpp\_int

```

namespace boost{ namespace multiprecision{

typedef unspecified-type limb_type;

enum cpp_integer_type    { signed_magnitude, unsigned_magnitude };
enum cpp_int_check_type  { checked, unchecked };

template <unsigned MinDigits = 0,
          unsigned MaxDits = 0,
          cpp_integer_type SignType = signed_magnitude,
          cpp_int_check_type Checked = unchecked,
          class Allocator = std::allocator<limb_type> >
class cpp_int_backend;
//
// Expression templates default to et_off if there is no allocator:
//
template <unsigned MinDigits, unsigned MaxDigits, cpp_integer_type Sign-
Type, cpp_int_check_type Checked>
struct expression_template_default<cpp_int_backend<MinDigits, MaxDigits, SignType, Checked, void> >
{ static const expression_template_option value = et_off; };

typedef number<cpp_int_backend<> >          cpp_int;    // arbitrary precision integer
typedef rational_adaptor<cpp_int_backend<> >  cpp_rational_backend;
typedef number<cpp_rational_backend>         cpp_rational; // arbitrary precision rational
number

// Fixed precision unsigned types:
typedef number<cpp_int_backend<128, 128, unsigned_magnitude, unchecked, void> >  uint128_t;
typedef number<cpp_int_backend<256, 256, unsigned_magnitude, unchecked, void> >  uint256_t;
typedef number<cpp_int_backend<512, 512, unsigned_magnitude, unchecked, void> >  uint512_t;
typedef number<cpp_int_backend<1024, 1024, unsigned_magnitude, unchecked, void> > uint1024_t;

// Fixed precision signed types:
typedef number<cpp_int_backend<128, 128, signed_magnitude, unchecked, void> >      int128_t;
typedef number<cpp_int_backend<256, 256, signed_magnitude, unchecked, void> >      int256_t;
typedef number<cpp_int_backend<512, 512, signed_magnitude, unchecked, void> >      int512_t;
typedef number<cpp_int_backend<1024, 1024, signed_magnitude, unchecked, void> >      int1024_t;

// Over again, but with checking enabled this time:
typedef number<cpp_int_backend<0, 0, signed_magnitude, checked> >                  checked_cpp_int;
typedef rational_adaptor<cpp_int_backend<0, 0, signed_magnitude, checked> >        checked_cpp_ra-
tional_backend;
typedef number<cpp_rational_backend>                                                checked_cpp_ra-
tional;

// Checked fixed precision unsigned types:
typedef number<cpp_int_backend<128, 128, unsigned_magnitude, checked, void> >      ↵
checked_uint128_t;
typedef number<cpp_int_backend<256, 256, unsigned_magnitude, checked, void> >      ↵
checked_uint256_t;
typedef number<cpp_int_backend<512, 512, unsigned_magnitude, checked, void> >      ↵
checked_uint512_t;
typedef number<cpp_int_backend<1024, 1024, unsigned_magnitude, checked, void> >      ↵
checked_uint1024_t;

// Fixed precision signed types:
typedef number<cpp_int_backend<128, 128, signed_magnitude, checked, void> >          checked_int128_t;

```

```
typedef number<cpp_int_backend<256, 256, signed_magnitude, checked, void> > checked_int256_t;
typedef number<cpp_int_backend<512, 512, signed_magnitude, checked, void> > checked_int512_t;
typedef number<cpp_int_backend<1024, 1024, signed_magnitude, checked, void> > checked_int1024_t;

}} // namespaces
```

Class template `cpp_int_backend` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The template arguments are:

MinBits	Determines the number of Bits to store directly within the object before resorting to dynamic memory allocation. When zero, this field is determined automatically based on how many bits can be stored in union with the dynamic storage header: setting a larger value may improve performance as larger integer values will be stored internally before memory allocation is required.
MaxBits	Determines the maximum number of bits to be stored in the type: resulting in a fixed precision type. When this value is the same as MinBits, then the Allocator parameter is ignored, as no dynamic memory allocation will ever be performed: in this situation the Allocator parameter should be set to type <code>void</code> . Note that this parameter should not be used simply to prevent large memory allocations, not only is that role better performed by the allocator, but fixed precision integers have a tendency to allocate all of MaxBits of storage more often than one would expect.
SignType	Determines whether the resulting type is signed or not. Note that for <a href="#">arbitrary precision</a> types this parameter must be <code>signed_magnitude</code> . For fixed precision types then this type may be either <code>signed_magnitude</code> or <code>unsigned_magnitude</code> .
Checked	This parameter has two values: <code>checked</code> or <code>unchecked</code> . See the <a href="#">tutorial</a> for more information.
Allocator	The allocator to use for dynamic memory allocation, or type <code>void</code> if <code>MaxBits == MinBits</code> .

The type of `number_category<cpp_int<Args...> >::type` is `mpl::int_<number_kind_integer>`.

More information on this type can be found in the [tutorial](#).

## gmp\_int

```
namespace boost{ namespace multiprecision{

class gmp_int;

typedef number<gmp_int > mpz_int;

}} // namespaces
```

Class template `gmp_int` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The type of `number_category<cpp_int<Args...> >::type` is `mpl::int_<number_kind_integer>`.

More information on this type can be found in the [tutorial](#).

## tom\_int

```
namespace boost{ namespace multiprecision{  
  
    class tommath_int;  
  
    typedef number<tommath_int >          tom_int;  
  
}} // namespaces
```

Class template `tommath_int` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The type of `number_category<cpp_int<Args...> >::type` is `mpl::int_<number_kind_integer>`.

More information on this type can be found in the [tutorial](#).

## gmp\_float

```
namespace boost{ namespace multiprecision{  
  
    template <unsigned Digits10>  
    class gmp_float;  
  
    typedef number<gmp_float<50> >      mpf_float_50;  
    typedef number<gmp_float<100> >     mpf_float_100;  
    typedef number<gmp_float<500> >     mpf_float_500;  
    typedef number<gmp_float<1000> >    mpf_float_1000;  
    typedef number<gmp_float<0> >       mpf_float;  
  
}} // namespaces
```

Class template `gmp_float` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The class takes a single template parameter - `Digits10` - which is the number of decimal digits precision the type should support. When this parameter is zero, then the precision can be set at runtime via `number::default_precision` and `number::precision`. Note that this type does not in any way change the GMP library's global state (for example it does not change the default precision of the `mpf_t` data type), therefore you can safely mix this type with existing code that uses GMP, and also mix `gmp_floats` of differing precision.

The type of `number_category<cpp_int<Args...> >::type` is `mpl::int_<number_kind_floating_point>`.

More information on this type can be found in the [tutorial](#).

## mpfr\_float\_backend

```
namespace boost{ namespace multiprecision{

template <unsigned Digits10>
class mpfr_float_backend;

typedef number<mpfr_float_backend<50> >      mpfr_float_50;
typedef number<mpfr_float_backend<100> >     mpfr_float_100;
typedef number<mpfr_float_backend<500> >     mpfr_float_500;
typedef number<mpfr_float_backend<1000> >    mpfr_float_1000;
typedef number<mpfr_float_backend<0> >      mpfr_float;

}} // namespaces
```

Class template `mpfr_float_backend` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The class takes a single template parameter - `Digits10` - which is the number of decimal digits precision the type should support. When this parameter is zero, then the precision can be set at runtime via `number::default_precision` and `number::precision`. Note that this type does not in any way change the GMP or MPFR library's global state (for example it does not change the default precision of the `mpfr_t` data type), therefore you can safely mix this type with existing code that uses GMP or MPFR, and also mix `mpfr_float_backends` of differing precision.

The type of `number_category<cpp_int<Args...> >::type` is `mpl::int_<number_kind_floating_point>`.

More information on this type can be found in the [tutorial](#).

## cpp\_bin\_float

```
namespace boost{ namespace multiprecision{

enum digit_base_type
{
    digit_base_2 = 2,
    digit_base_10 = 10
};

template <unsigned Digits, digit_base_type base = digit_base_10, class Allocator = void, class Exponent = int, ExponentMin = 0, ExponentMax = 0>
class cpp_bin_float;

typedef number<cpp_bin_float<50> > cpp_bin_float_50;
typedef number<cpp_bin_float<100> > cpp_bin_float_100;

typedef number<backends::cpp_bin_float<24, backends::digit_base_2, void, boost::int16_t, -126, 127>, et_off>      cpp_bin_float_single;
typedef number<backends::cpp_bin_float<53, backends::digit_base_2, void, boost::int16_t, -1022, 1023>, et_off>      cpp_bin_float_double;
typedef number<backends::cpp_bin_float<64, backends::digit_base_2, void, boost::int16_t, -16382, 16383>, et_off>      cpp_bin_float_double_extended;
typedef number<backends::cpp_bin_float<113, backends::digit_base_2, void, boost::int16_t, -16382, 16383>, et_off>      cpp_bin_float_quad;

}} // namespaces
```

Class template `cpp_bin_float` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The class takes six template parameters:

Digits	The number of digits precision the type should support. This is normally expressed as base-10 digits, but that can be changed via the second template parameter.
base	An enumerated value (either <code>digit_base_10</code> or <code>digit_base_2</code> ) that indicates whether <code>Digits</code> is base-10 or base-2
Allocator	The allocator used: defaults to type <code>void</code> , meaning all storage is within the class, and no dynamic allocation is performed, but can be set to a standard library allocator if dynamic allocation makes more sense.
Exponent	A signed integer type to use as the type of the exponent - defaults to <code>int</code> .
ExponentMin	The smallest (most negative) permitted exponent, defaults to zero, meaning "define as small as possible given the limitations of the type and our internal requirements".
ExponentMax	The largest (most positive) permitted exponent, defaults to zero, meaning "define as large as possible given the limitations of the type and our internal requirements".

The type of `number_category<cpp_bin_float<Args...>::type` is `mpl::int_<number_kind_floating_point>`.

More information on this type can be found in the [tutorial](#).

## Implementation Notes

Internally, an N-bit `cpp_bin_float` is represented as an N-bit unsigned integer along with an exponent and a sign. The integer part is normalized so that it's most significant bit is always 1. The decimal point is assumed to be directly after the most significant bit of the integer part. The special values zero, infinity and NaN all have the integer part set to zero, and the exponent to one of 3 special values above the maximum permitted exponent.

Multiplication is trivial: multiply the two N-bit integer mantissa's to obtain a 2N-bit number, then round and adjust the sign and exponent.

Addition and subtraction proceed similarly - if the exponents are such that there is overlap between the two values, then left shift the larger value to produce a number with between N and 2N bits, then perform integer addition or subtraction, round, and adjust the exponent.

Division proceeds as follows: first scale the numerator by some power of 2 so that integer division will produce either an N-bit or N+1 bit result plus a remainder. If we get an N bit result then the size of twice the remainder compared to the denominator gives us the rounding direction. Otherwise we have one extra bit in the result which we can use to determine rounding (in this case ties occur only if the remainder is zero and the extra bit is a 1).

Square root uses integer square root in a manner analogous to division.

Decimal string to binary conversion proceeds as follows: first parse the digits to produce an integer multiplied by a decimal exponent. Note that we stop parsing digits once we have parsed as many as can possibly effect the result - this stops the integer part growing too large when there are a very large number of input digits provided. At this stage if the decimal exponent is positive then the result is an integer and we can in principle simply multiply by  $10^N$  to get an exact integer result. In practice however, that could produce some very large integers. We also need to be able to divide by  $10^N$  in the event that the exponent is negative. Therefore calculation of the  $10^N$  values plus the multiplication or division are performed using limited precision integer arithmetic, plus an exponent, and a track of the accumulated error. At the end of the calculation we will either be able to round unambiguously, or the error will be such that we can't tell which way to round. In the latter case we simply up the precision and try again until we have an unambiguously rounded result.

Binary to decimal conversion proceeds very similarly to the above, our aim is to calculate `mantissa * 2shift * 10E` where `E` is the decimal exponent and `shift` is calculated so that the result is an N bit integer assuming we want N digits printed in the result. As before we use limited precision arithmetic to calculate the result and up the precision as necessary until the result is unambiguously correctly rounded. In addition our initial calculation of the decimal exponent may be out by 1, so we have to correct that and loop as well in the that case.

## cpp\_dec\_float

```
namespace boost{ namespace multiprecision{  
  
    template <unsigned Digits10, class ExponentType = boost::int32_t, class Allocator = void>  
    class cpp_dec_float;  
  
    typedef number<cpp_dec_float<50> > cpp_dec_float_50;  
    typedef number<cpp_dec_float<100> > cpp_dec_float_100;  
  
}} // namespaces
```

Class template `cpp_dec_float` fulfils all of the requirements for a [Backend](#) type. Its members and non-member functions are deliberately not documented: these are considered implementation details that are subject to change.

The class takes three template parameters:

Digits10	The number of decimal digits precision the type should support. Note that this type does not normally perform any dynamic memory allocation, and as a result the <code>Digits10</code> template argument should not be set too high or the class's size will grow unreasonably large.
ExponentType	A signed integer type that represents the exponent of the number
Allocator	The allocator used: defaults to type <code>void</code> , meaning all storage is within the class, and no dynamic allocation is performed, but can be set to a standard library allocator if dynamic allocation makes more sense.

The type of `number_category<cpp_dec_float<Args...>::type` is `mpl::int_<number_kind_floating_point>`.

More information on this type can be found in the [tutorial](#).

## Internal Support Code

There are some traits classes which authors of new backends should be aware of:

```
namespace boost{ namespace multiprecision{ namespace detail{  
  
    template<typename From, typename To>  
    struct is_explicitly_convertible;  
  
}}}
```

Inherits from `boost::integral_constant<bool, true>` if type `From` has an explicit conversion from `To`.

For compilers that support C++11 SFINAE-expressions this trait should "just work". Otherwise it inherits from `boost::is_convertible<From, To>::type`, and will need to be specialised for Backends that have constructors marked as `explicit`.

```
template <class From, class To>  
struct is_lossy_conversion  
{  
    static const bool value = see below;  
};
```

Member value is true if the conversion from `From` to `To` would result in a loss of precision, and false otherwise.

The default version of this trait simply checks whether the *kind* of conversion (for example from a floating point to an integer type) is inherently lossy. Note that if either of the types `From` or `To` are of an unknown number category (because `number_category` is not specialised for that type) then this trait will be true.

```
template<typename From, typename To>
struct is_restricted_conversion
{
    static const bool value = see below;
};
```

Member value is true if From is only explicitly convertible to To and not implicitly convertible, or if `is_lossy_conversion<From, To>::value` is true. Otherwise false.

Note that while this trait is the ultimate arbiter of which constructors are marked as `explicit` in class `number`, authors of backend types should generally specialise one of the traits above, rather than this one directly.

```
template <class T>
is_signed_number;
template <class T>
is_unsigned_number;
```

These two traits inherit from either `mpl::true_` or `mpl::false_`, by default types are assumed to be signed unless `is_unsigned_number` is specialized for that type.

## Backend Requirements

The requirements on the Backend template argument to `number` are split up into sections: compulsory and optional.

Compulsory requirements have no default implementation in the library, therefore if the feature they implement is to be supported at all, then they must be implemented by the backend.

Optional requirements have default implementations that are called if the backend doesn't provide it's own. Typically the backend will implement these to improve performance.

In the following tables, type `B` is the Backend template argument to `number`, `b` and `b2` are a variables of type `B`, `cb`, `cb2` and `cb3` are constant variables of type `const B`, `rb` is a variable of type `B&&`, `a` and `a2` are variables of Arithmetic type, `s` is a variable of type `const char*`, `ui` is a variable of type `unsigned`, `bb` is a variable of type `bool`, `pa` is a variable of type pointer-to-arithmetic-type, `exp` is a variable of type `B::exp_type`, `pexp` is a variable of type `B::exp_type*`, `i` is a variable of type `int`, `pi` pointer to a variable of type `int`, `B2` is another type that meets these requirements, `b2` is a variable of type `B2`, `ss` is variable of type `std::streamsize` and `ff` is a variable of type `std::ios_base::fmtflags`.

**Table 8. Compulsory Requirements on the Backend type.**

Expression	Return Type	Comments	Throws
<code>B::signed_types</code>	<code>mpl::list&lt;type-list&gt;</code>	A list of signed integral types that can be assigned to type B. The types shall be listed in order of size, smallest first, and shall terminate in the type that is <code>std::intmax_t</code> .	
<code>B::unsigned_types</code>	<code>mpl::list&lt;type-list&gt;</code>	A list of unsigned integral types that can be assigned to type B. The types shall be listed in order of size, smallest first, and shall terminate in the type that is <code>std::uintmax_t</code> .	
<code>B::float_types</code>	<code>mpl::list&lt;type-list&gt;</code>	A list of floating-point types that can be assigned to type B. The types shall be listed in order of size, smallest first, and shall terminate in type <code>long double</code> .	
<code>B::exponent_type</code>	A signed integral type.	The type of the exponent of type B. This type is required only for floating point types.	
<code>B()</code>		Default constructor.	
<code>B(cb)</code>		Copy Constructor.	
<code>b = b</code>	<code>B&amp;</code>	Assignment operator.	
<code>b = a</code>	<code>B&amp;</code>	Assignment from an Arithmetic type. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> .	
<code>b = s</code>	<code>B&amp;</code>	Assignment from a string.	Throws a <code>std::runtime_error</code> if the string could not be interpreted as a valid number.
<code>b.swap(b)</code>	<code>void</code>	Swaps the contents of its arguments.	<code>noexcept</code>
<code>cb.str(ss, ff)</code>	<code>std::string</code>	Returns the string representation of b with <code>ss</code> digits and formatted according to the flags set in <code>ff</code> . If <code>ss</code> is zero, then returns as many digits as are required to reconstruct the original value.	



Expression	Return Type	Comments	Throws
<code>b.negate()</code>	<code>void</code>	Negates <code>b</code> .	
<code>cb.compare(cb2)</code>	<code>int</code>	Compares <code>cb</code> and <code>cb2</code> , returns a value less than zero if <code>cb &lt; cb2</code> , a value greater than zero if <code>cb &gt; cb2</code> and zero if <code>cb == cb2</code> .	<code>noexcept</code>
<code>cb.compare(a)</code>	<code>int</code>	Compares <code>cb</code> and <code>a</code> , returns a value less than zero if <code>cb &lt; a</code> , a value greater than zero if <code>cb &gt; a</code> and zero if <code>cb == a</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> .	
<code>eval_add(b, cb)</code>	<code>void</code>	Adds <code>cb</code> to <code>b</code> .	
<code>eval_subtract(b, cb)</code>	<code>void</code>	Subtracts <code>cb</code> from <code>b</code> .	
<code>eval_multiply(b, cb)</code>	<code>void</code>	Multiplies <code>b</code> by <code>cb</code> .	
<code>eval_divide(b, cb)</code>	<code>void</code>	Divides <code>b</code> by <code>cb</code> .	<code>std::overflow_error</code> if <code>cb</code> has the value zero, and <code>std::numeric_limits&lt;number&lt;B&gt;&gt;::has_infinity == false</code>
<code>eval_modulus(b, cb)</code>	<code>void</code>	Computes <code>b %= cb</code> , only required when <code>B</code> is an integer type.	<code>std::overflow_error</code> if <code>cb</code> has the value zero.
<code>eval_bitwise_and(b, cb)</code>	<code>void</code>	Computes <code>b &amp;= cb</code> , only required when <code>B</code> is an integer type.	
<code>eval_bitwise_or(b, cb)</code>	<code>void</code>	Computes <code>b  = cb</code> , only required when <code>B</code> is an integer type.	
<code>eval_bitwise_xor(b, cb)</code>	<code>void</code>	Computes <code>b ^= cb</code> , only required when <code>B</code> is an integer type.	
<code>eval_complement(b, cb)</code>	<code>void</code>	Computes the ones-complement of <code>cb</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type.	
<code>eval_left_shift(b, ui)</code>	<code>void</code>	Computes <code>b &lt;&lt;= ui</code> , only required when <code>B</code> is an integer type.	

Expression	Return Type	Comments	Throws
<code>eval_right_shift(b, ui)</code>	<code>void</code>	Computes $b \gg= ui$ , only required when <code>B</code> is an integer type.	
<code>eval_convert_to(pa, cb)</code>	<code>void</code>	Converts <code>cb</code> to the type of <code>*pa</code> and store the result in <code>*pa</code> . Type <code>B</code> shall support conversion to at least types <code>std::intmax_t</code> , <code>std::uintmax_t</code> and <code>long long</code> . Conversion to other arithmetic types can then be synthesised using other operations. Conversions to other types are entirely optional.	
<code>eval_frexp(b, cb, pexp)</code>	<code>void</code>	Stores values in <code>b</code> and <code>*pexp</code> such that the value of <code>cb</code> is $b * 2^{pexp}$ , only required when <code>B</code> is a floating-point type.	
<code>eval_ldexp(b, cb, exp)</code>	<code>void</code>	Stores a value in <code>b</code> that is $cb * 2^{exp}$ , only required when <code>B</code> is a floating-point type.	
<code>eval_frexp(b, cb, pi)</code>	<code>void</code>	Stores values in <code>b</code> and <code>*pi</code> such that the value of <code>cb</code> is $b * 2^{pi}$ , only required when <code>B</code> is a floating-point type.	<code>std::runtime_error</code> if the exponent of <code>cb</code> is too large to be stored in an <code>int</code> .
<code>eval_ldexp(b, cb, i)</code>	<code>void</code>	Stores a value in <code>b</code> that is $cb * 2^i$ , only required when <code>B</code> is a floating-point type.	
<code>eval_floor(b, cb)</code>	<code>void</code>	Stores the floor of <code>cb</code> in <code>b</code> , only required when <code>B</code> is a floating-point type.	
<code>eval_ceil(b, cb)</code>	<code>void</code>	Stores the ceiling of <code>cb</code> in <code>b</code> , only required when <code>B</code> is a floating-point type.	
<code>eval_sqrt(b, cb)</code>	<code>void</code>	Stores the square root of <code>cb</code> in <code>b</code> , only required when <code>B</code> is a floating-point type.	

Expression	Return Type	Comments	Throws
<code>boost::multiprecision::number_category&lt;B&gt;::type</code>	<code>mpl::int_&lt;N&gt;</code>	N is one of the values <code>number_kind_integer</code> , <code>number_kind_floating_point</code> , <code>number_kind_rational</code> or <code>number_kind_fixed_point</code> . Defaults to <code>number_kind_floating_point</code> .	

**Table 9. Optional Requirements on the Backend Type**

Expression	Returns	Comments	Throws
<i>Construct and assign:</i>			
<code>B(rb)</code>	<code>B</code>	Move constructor. Afterwards variable <code>rb</code> shall be in sane state, albeit with unspecified value. Only destruction and assignment to the moved-from variable <code>rb</code> need be supported after the operation.	<code>noexcept</code>
<code>b = rb</code>	<code>B&amp;</code>	Move-assign. Afterwards variable <code>rb</code> shall be in sane state, albeit with unspecified value. Only destruction and assignment to the moved-from variable <code>rb</code> need be supported after the operation.	<code>noexcept</code>
<code>B(a)</code>	<code>B</code>	Direct construction from an arithmetic type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, this operation is simulated using default-construction followed by assignment.	
<code>B(b2)</code>	<code>B</code>	Copy constructor from a different back-end type. When not provided, a generic interconversion routine is used. This constructor may be <code>explicit</code> if the corresponding frontend constructor should also be <code>explicit</code> .	
<code>b = b2</code>	<code>b&amp;</code>	Assignment operator from a different back-end type. When not provided, a generic interconversion routine is used.	
<code>assign_components(b, a, a)</code>	<code>void</code>	Assigns to <code>b</code> the two components in the following arguments. Only applies to rational and complex number types. When not provided, arithmetic operations are used to synthesise the result from the two values.	

Expression	Returns	Comments	Throws
<code>assign_components(b, b2, b2)</code>	<code>void</code>	Assigns to <code>b</code> the two components in the following arguments. Only applies to rational and complex number types. When not provided, arithmetic operations are used to synthesise the result from the two values.	
<i>Comparisons:</i>			
<code>eval_eq(cb, cb2)</code>	<code>bool</code>	Returns <code>true</code> if <code>cb</code> and <code>cb2</code> are equal in value. When not provided, the default implementation returns <code>cb.compare(cb2) == 0</code> .	<code>noexcept</code>
<code>eval_eq(cb, a)</code>	<code>bool</code>	Returns <code>true</code> if <code>cb</code> and <code>a</code> are equal in value. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, return the equivalent of <code>eval_eq(cb, B(a))</code> .	
<code>eval_eq(a, cb)</code>	<code>bool</code>	Returns <code>true</code> if <code>cb</code> and <code>a</code> are equal in value. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version returns <code>eval_eq(cb, a)</code> .	
<code>eval_lt(cb, cb2)</code>	<code>bool</code>	Returns <code>true</code> if <code>cb</code> is less than <code>cb2</code> in value. When not provided, the default implementation returns <code>cb.compare(cb2) &lt; 0</code> .	<code>noexcept</code>
<code>eval_lt(cb, a)</code>	<code>bool</code>	Returns <code>true</code> if <code>cb</code> is less than <code>a</code> in value. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default implementation returns <code>eval_lt(cb, B(a))</code> .	

Expression	Returns	Comments	Throws
<code>eval_lt(a, cb)</code>	bool	Returns true if a is less than cb in value. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default implementation returns <code>eval_gt(cb, a)</code> .	
<code>eval_gt(cb, cb2)</code>	bool	Returns true if cb is greater than cb2 in value. When not provided, the default implementation returns <code>cb.compare(cb2) &gt; 0</code> .	noexcept
<code>eval_gt(cb, a)</code>	bool	Returns true if cb is greater than a in value. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default implementation returns <code>eval_gt(cb, B(a))</code> .	
<code>eval_gt(a, cb)</code>	bool	Returns true if a is greater than cb in value. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default implementation returns <code>eval_lt(cb, a)</code> .	
<code>eval_is_zero(cb)</code>	bool	Returns true if cb is zero, otherwise false. The default version of this function returns <code>cb.compare(ui_type(0)) == 0</code> , where <code>ui_type</code> is <code>typename mpl::front&lt;typename B::unsigned_types&gt;::type</code> .	

Expression	Returns	Comments	Throws
<code>eval_get_sign(cb)</code>	<code>int</code>	Returns a value < zero if <code>cb</code> is negative, a value > zero if <code>cb</code> is positive, and zero if <code>cb</code> is zero. The default version of this function returns <code>cb.compare(ui_type(0))</code> , where <code>ui_type</code> is <code>ui_type is type_name mpl::front&lt;typename B::unsigned_types&gt;::type</code> .	
<i>Basic arithmetic:</i>			
<code>eval_add(b, a)</code>	<code>void</code>	Adds <code>a</code> to <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_add(b, B(a))</code>	
<code>eval_add(b, cb, cb2)</code>	<code>void</code>	Add <code>cb</code> to <code>cb2</code> and stores the result in <code>b</code> . When not provided, does the equivalent of <code>b = cb; eval_add(b, cb2)</code> .	
<code>eval_add(b, cb, a)</code>	<code>void</code>	Add <code>cb</code> to <code>a</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_add(b, cb, B(a))</code> .	
<code>eval_add(b, a, cb)</code>	<code>void</code>	Add <code>a</code> to <code>cb</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_add(b, cb, a)</code> .	
<code>eval_subtract(b, a)</code>	<code>void</code>	Subtracts <code>a</code> from <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_subtract(b, B(a))</code>	

Expression	Returns	Comments	Throws
<code>eval_subtract(b, cb, cb2)</code>	void	Subtracts <code>cb2</code> from <code>cb</code> and stores the result in <code>b</code> . When not provided, does the equivalent of <code>b = cb; eval_subtract(b, cb2)</code> .	
<code>eval_subtract(b, cb, a)</code>	void	Subtracts <code>a</code> from <code>cb</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_subtract(b, cb, B(a))</code> .	
<code>eval_subtract(b, a, cb)</code>	void	Subtracts <code>cb</code> from <code>a</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_subtract(b, cb, a); b.negate();</code> .	
<code>eval_multiply(b, a)</code>	void	Multiplies <code>b</code> by <code>a</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_multiply(b, B(a))</code> .	
<code>eval_multiply(b, cb, cb2)</code>	void	Multiplies <code>cb</code> by <code>cb2</code> and stores the result in <code>b</code> . When not provided, does the equivalent of <code>b = cb; eval_multiply(b, cb2)</code> .	
<code>eval_multiply(b, cb, a)</code>	void	Multiplies <code>cb</code> by <code>a</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_multiply(b, cb, B(a))</code> .	



Expression	Returns	Comments	Throws
<code>eval_multiply(b, a, cb)</code>	void	Multiplies a by cb and stores the result in b. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_multiply(b, cb, a)</code> .	
<code>eval_multiply_add(b, cb, cb2)</code>	void	Multiplies cb by cb2 and adds the result to b. When not provided does the equivalent of creating a temporary B t and <code>eval_multiply(t, cb, cb2)</code> followed by <code>eval_add(b, t)</code> .	
<code>eval_multiply_add(b, cb, a)</code>	void	Multiplies a by cb and adds the result to b. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided does the equivalent of creating a temporary B t and <code>eval_multiply(t, cb, a)</code> followed by <code>eval_add(b, t)</code> .	
<code>eval_multiply_add(b, a, cb)</code>	void	Multiplies a by cb and adds the result to b. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided does the equivalent of <code>eval_multiply_add(b, cb, a)</code> .	
<code>eval_multiply_subtract(b, cb, cb2)</code>	void	Multiplies cb by cb2 and subtracts the result from b. When not provided does the equivalent of creating a temporary B t and <code>eval_multiply(t, cb, cb2)</code> followed by <code>eval_subtract(b, t)</code> .	

Expression	Returns	Comments	Throws
<code>eval_multiply_subtract(b, cb, a)</code>	void	Multiplies a by cb and subtracts the result from b. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided does the equivalent of creating a temporary B t and <code>eval_multiply(t, cb, a)</code> followed by <code>eval_subtract(b, t)</code> .	
<code>eval_multiply_subtract(b, a, cb)</code>	void	Multiplies a by cb and subtracts the result from b. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided does the equivalent of <code>eval_multiply_subtract(b, cb, a)</code> .	
<code>eval_multiply_add(b, cb, cb2, cb3)</code>	void	Multiplies cb by cb2 and adds the result to cb3 storing the result in b. When not provided does the equivalent of <code>eval_multiply(b, cb, cb2)</code> followed by <code>eval_add(b, cb3)</code> . For brevity, only a version showing all arguments of type B is shown here, but you can replace up to any 2 of cb, cb2 and cb3 with any type listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> .	
<code>eval_multiply_subtract(b, cb, cb2, cb3)</code>	void	Multiplies cb by cb2 and subtracts from the result cb3 storing the result in b. When not provided does the equivalent of <code>eval_multiply(b, cb, cb2)</code> followed by <code>eval_subtract(b, cb3)</code> . For brevity, only a version showing all arguments of type B is shown here, but you can replace up to any 2 of cb, cb2 and cb3 with any type listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> .	

Expression	Returns	Comments	Throws
<code>eval_divide(b, a)</code>	<code>void</code>	Divides <code>b</code> by <code>a</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_divide(b, B(a))</code>	<code>std::overflow_error</code> if <code>a</code> has the value zero, and <code>std::numeric_limits&lt;number&lt;B&gt; &gt;::has_infinity == false</code>
<code>eval_divide(b, cb, cb2)</code>	<code>void</code>	Divides <code>cb</code> by <code>cb2</code> and stores the result in <code>b</code> . When not provided, does the equivalent of <code>b = cb; eval_divide(b, cb2)</code> .	<code>std::overflow_error</code> if <code>cb2</code> has the value zero, and <code>std::numeric_limits&lt;number&lt;B&gt; &gt;::has_infinity == false</code>
<code>eval_divide(b, cb, a)</code>	<code>void</code>	Divides <code>cb</code> by <code>a</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_divide(b, cb, B(a))</code> .	<code>std::overflow_error</code> if <code>a</code> has the value zero, and <code>std::numeric_limits&lt;number&lt;B&gt; &gt;::has_infinity == false</code>
<code>eval_divide(b, a, cb)</code>	<code>void</code>	Divides <code>a</code> by <code>cb</code> and stores the result in <code>b</code> . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_divide(b, B(a), cb)</code> .	<code>std::overflow_error</code> if <code>cb</code> has the value zero, and <code>std::numeric_limits&lt;number&lt;B&gt; &gt;::has_infinity == false</code>
<code>eval_increment(b)</code>	<code>void</code>	Increments the value of <code>b</code> by one. When not provided, does the equivalent of <code>eval_add(b, static_cast&lt;ui_type&gt;(1u))</code> . Where <code>ui_type</code> is <code>typename mpl::front&lt;typename B : : unsigned_types&gt;::type</code> .	
<code>eval_decrement(b)</code>	<code>void</code>	Decrements the value of <code>b</code> by one. When not provided, does the equivalent of <code>eval_subtract(b, static_cast&lt;ui_type&gt;(1u))</code> . Where <code>ui_type</code> is <code>typename mpl::front&lt;typename B : : unsigned_types&gt;::type</code> .	

Expression	Returns	Comments	Throws
<i>Integer specific operations:</i>			
<code>eval_modulus(b, a)</code>	void	Computes $b \mathrel{%=} cb$ , only required when B is an integer type. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_modulus(b, B(a))</code>	<code>std::overflow_error</code> if a has the value zero.
<code>eval_modulus(b, cb, cb2)</code>	void	Computes $cb \mathrel{\%} cb2$ and stores the result in b, only required when B is an integer type. When not provided, does the equivalent of $b = cb; \text{eval\_modulus}(b, cb2)$ .	<code>std::overflow_error</code> if a has the value zero.
<code>eval_modulus(b, cb, a)</code>	void	Computes $cb \mathrel{\%} a$ and stores the result in b, only required when B is an integer type. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_modulus(b, cb, B(a))</code> .	<code>std::overflow_error</code> if a has the value zero.
<code>eval_modulus(b, a, cb)</code>	void	Computes $cb \mathrel{\%} a$ and stores the result in b, only required when B is an integer type. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_modulus(b, B(a), cb)</code> .	<code>std::overflow_error</code> if a has the value zero.
<code>eval_bitwise_and(b, a)</code>	void	Computes $b \mathrel{\&=} cb$ , only required when B is an integer type. The type of a shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_bitwise_and(b, B(a))</code>	

Expression	Returns	Comments	Throws
<code>eval_bitwise_and(b, cb, cb2)</code>	void	Computes <code>cb &amp; cb2</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. When not provided, does the equivalent of <code>b = cb; eval_bitwise_and(b, cb2)</code> .	
<code>eval_bitwise_and(b, cb, a)</code>	void	Computes <code>cb &amp; a</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_and(b, cb, B(a))</code> .	
<code>eval_bitwise_and(b, a, cb)</code>	void	Computes <code>cb &amp; a</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_and(b, cb, a)</code> .	
<code>eval_bitwise_or(b, a)</code>	void	Computes <code>b  = cb</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_bitwise_or(b, B(a))</code> .	
<code>eval_bitwise_or(b, cb, cb2)</code>	void	Computes <code>cb   cb2</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. When not provided, does the equivalent of <code>b = cb; eval_bitwise_or(b, cb2)</code> .	

Expression	Returns	Comments	Throws
<code>eval_bitwise_or(b, cb, a)</code>	void	Computes $cb \mid a$ and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_or(b, cb, B(a))</code> .	
<code>eval_bitwise_or(b, a, cb)</code>	void	Computes $cb \mid a$ and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_or(b, cb, a)</code> .	
<code>eval_bitwise_xor(b, a)</code>	void	Computes $b \wedge cb$ , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, the default version calls <code>eval_bitwise_xor(b, B(a))</code> .	
<code>eval_bitwise_xor(b, cb, cb2)</code>	void	Computes $cb \wedge cb2$ and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. When not provided, does the equivalent of <code>b = cb; eval_bitwise_xor(b, cb2)</code> .	
<code>eval_bitwise_xor(b, cb, a)</code>	void	Computes $cb \wedge a$ and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_xor(b, cb, B(a))</code> .	

Expression	Returns	Comments	Throws
<code>eval_bitwise_xor(b, a, cb)</code>	void	Computes $a \wedge cb$ and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . When not provided, does the equivalent of <code>eval_bitwise_xor(b, cb, a)</code> .	
<code>eval_left_shift(b, cb, ui)</code>	void	Computes <code>cb &lt;&lt; ui</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. When not provided, does the equivalent of <code>b = cb; eval_left_shift(b, a);</code> .	
<code>eval_right_shift(b, cb, ui)</code>	void	Computes <code>cb &gt;&gt; ui</code> and stores the result in <code>b</code> , only required when <code>B</code> is an integer type. When not provided, does the equivalent of <code>b = cb; eval_right_shift(b, a);</code> .	
<code>eval_qr(cb, cb2, b, b2)</code>	void	Sets <code>b</code> to the result of <code>cb / cb2</code> and <code>b2</code> to the result of <code>cb % cb2</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	<code>std::overflow_error</code> if <code>a</code> has the value zero.
<code>eval_integer_modulus(cb, ui)</code>	unsigned	Returns the result of <code>cb % ui</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	<code>std::overflow_error</code> if <code>a</code> has the value zero.
<code>eval_lsb(cb)</code>	unsigned	Returns the index of the least significant bit that is set. Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_msb(cb)</code>	unsigned	Returns the index of the most significant bit that is set. Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	

Expression	Returns	Comments	Throws
<code>eval_bit_test(cb, ui)</code>	<code>bool</code>	Returns true if <code>cb</code> has bit <code>ui</code> set. Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_bit_set(b, ui)</code>	<code>void</code>	Sets the bit at index <code>ui</code> in <code>b</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_bit_unset(b, ui)</code>	<code>void</code>	Unsets the bit at index <code>ui</code> in <code>b</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_bit_flip(b, ui)</code>	<code>void</code>	Flips the bit at index <code>ui</code> in <code>b</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_gcd(b, cb, cb2)</code>	<code>void</code>	Sets <code>b</code> to the greatest common divisor of <code>cb</code> and <code>cb2</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_lcm(b, cb, cb2)</code>	<code>void</code>	Sets <code>b</code> to the least common multiple of <code>cb</code> and <code>cb2</code> . Only required when <code>B</code> is an integer type. The default version of this function is synthesised from other operations above.	
<code>eval_gcd(b, cb, a)</code>	<code>void</code>	Sets <code>b</code> to the greatest common divisor of <code>cb</code> and <code>cb2</code> . Only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . The default version of this function calls <code>eval_gcd(b, cb, B(a))</code> .	



Expression	Returns	Comments	Throws
<code>eval_lcm(b, cb, a)</code>	void	Sets <code>b</code> to the least common multiple of <code>cb</code> and <code>cb2</code> . Only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . The default version of this function calls <code>eval_lcm(b, cb, B(a))</code> .	
<code>eval_gcd(b, a, cb)</code>	void	Sets <code>b</code> to the greatest common divisor of <code>cb</code> and <code>a</code> . Only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . The default version of this function calls <code>eval_gcd(b, cb, a)</code> .	
<code>eval_lcm(b, a, cb)</code>	void	Sets <code>b</code> to the least common multiple of <code>cb</code> and <code>a</code> . Only required when <code>B</code> is an integer type. The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> or <code>B::float_types</code> . The default version of this function calls <code>eval_lcm(b, cb, a)</code> .	
<code>eval_powm(b, cb, cb2, cb3)</code>	void	Sets <code>b</code> to the result of $(cb^{cb2})\%cb3$ . The default version of this function is synthesised from other operations above.	
<code>eval_powm(b, cb, cb2, a)</code>	void	Sets <code>b</code> to the result of $(cb^{cb2})\%a$ . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> . The default version of this function is synthesised from other operations above.	

Expression	Returns	Comments	Throws
<code>eval_powm(b, cb, a, cb2)</code>	void	Sets <code>b</code> to the result of $(cb^a)\%cb2$ . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> . The default version of this function is synthesised from other operations above.	
<code>eval_powm(b, cb, a, a2)</code>	void	Sets <code>b</code> to the result of $(cb^a)\%a2$ . The type of <code>a</code> shall be listed in one of the type lists <code>B::signed_types</code> , <code>B::unsigned_types</code> . The default version of this function is synthesised from other operations above.	
<code>eval_integer_sqrt(b, cb, b2)</code>	void	Sets <code>b</code> to the largest integer which when squared is less than <code>cb</code> , also sets <code>b2</code> to the remainder, ie to $cb - b^2$ . The default version of this function is synthesised from other operations above.	
<i>Sign manipulation:</i>			
<code>eval_abs(b, cb)</code>	void	Set <code>b</code> to the absolute value of <code>cb</code> . The default version of this functions assigns <code>cb</code> to <code>b</code> , and then calls <code>b.negate()</code> if <code>eval_get_sign(cb) &lt; 0</code> .	
<code>eval_fabs(b, cb)</code>	void	Set <code>b</code> to the absolute value of <code>cb</code> . The default version of this functions assigns <code>cb</code> to <code>b</code> , and then calls <code>b.negate()</code> if <code>eval_get_sign(cb) &lt; 0</code> .	
<i>Floating point functions:</i>			
<code>eval_fpclassify(cb)</code>	int	Returns one of the same values returned by <code>std::fpclassify</code> . Only required when <code>B</code> is an floating-point type. The default version of this function will only test for zero <code>cb</code> .	

Expression	Returns	Comments	Throws
<code>eval_trunc(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::trunc</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_round(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::round</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_exp(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::exp</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_log(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::log</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_log10(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::log10</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_sin(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::sin</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	

Expression	Returns	Comments	Throws
<code>eval_cos(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::cos</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_tan(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::exp</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_asin(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::asin</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_acos(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::acos</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_atan(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::atan</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_sinh(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::sinh</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	

Expression	Returns	Comments	Throws
<code>eval_cosh(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::cosh</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_tanh(b, cb)</code>	<code>void</code>	Performs the equivalent operation to <code>std::tanh</code> on argument <code>cb</code> and stores the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_fmod(b, cb, cb2)</code>	<code>void</code>	Performs the equivalent operation to <code>std::fmod</code> on arguments <code>cb</code> and <code>cb2</code> , and store the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_pow(b, cb, cb2)</code>	<code>void</code>	Performs the equivalent operation to <code>std::pow</code> on arguments <code>cb</code> and <code>cb2</code> , and store the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	
<code>eval_atan2(b, cb, cb2)</code>	<code>void</code>	Performs the equivalent operation to <code>std::atan</code> on arguments <code>cb</code> and <code>cb2</code> , and store the result in <code>b</code> . Only required when <code>B</code> is an floating-point type. The default version of this function is synthesised from other operations above.	

When the tables above place no *throws* requirements on an operation, then it is up to each type modelling this concept to decide when or whether throwing an exception is desirable. However, thrown exceptions should always either be the type, or inherit from the type `std::runtime_error`. For example, a floating point type might choose to throw `std::overflow_error` whenever the result of an operation would be infinite, and `std::underflow_error` whenever it would round to zero.



### Note

The non-member functions are all named with an "eval\_" prefix to avoid conflicts with template classes of the same name - in point of fact this naming convention shouldn't be necessary, but rather works around some compiler bugs.

## Header File Structure

**Table 10.** Top level headers

Header	Contains
cpp_int.hpp	The <code>cpp_int</code> backend type.
gmp.hpp	Defines all <a href="#">GMP</a> related backends.
millerrabin.hpp	Miller Rabin primality testing code.
number.hpp	Defines the <code>number</code> backend, is included by all the backend headers.
mpfr.hpp	Defines the <code>mpfr_float_backend</code> backend.
random.hpp	Defines code to interoperate with <code>Boost.Random</code> .
rational_adaptor.hpp	Defines the <code>rational_adaptor</code> backend.
cpp_dec_float.hpp	Defines the <code>cpp_dec_float</code> backend.
tommath.hpp	Defines the <code>tommath_int</code> backend.
concepts/number_archetypes.hpp	Defines a backend concept archetypes for testing use.

**Table 11. Implementation Headers**

Header	Contains
cpp_int/add.hpp	Add and subtract operators for <code>cpp_int_backend</code> .
cpp_int/bitwise.hpp	Bitwise operators for <code>cpp_int_backend</code> .
cpp_int/checked.hpp	Helper functions for checked arithmetic for <code>cpp_int_backend</code> .
cpp_int/comparison.hpp	Comparison operators for <code>cpp_int_backend</code> .
cpp_int/cpp_int_config.hpp	Basic setup and configuration for <code>cpp_int_backend</code> .
cpp_int/divide.hpp	Division and modulus operators for <code>cpp_int_backend</code> .
cpp_int/limits.hpp	<code>numeric_limits</code> support for <code>cpp_int_backend</code> .
cpp_int/misc.hpp	Miscellaneous operators for <code>cpp_int_backend</code> .
cpp_int/multiply.hpp	Multiply operators for <code>cpp_int_backend</code> .
detail/big_lanczos.hpp	Lanczos support for Boost.Math integration.
detail/default_ops.hpp	Default versions of the optional backend non-member functions.
detail/generic_interconvert.hpp	Generic interconversion routines.
detail/number_base.hpp	All the expression template code, metaprogramming, and operator overloads for <code>number</code> .
detail/no_et_ops.hpp	The non-expression template operators.
detail/functions/constants.hpp	Defines constants used by the floating point functions.
detail/functions/pow.hpp	Defines default versions of the power and exponential related floating point functions.
detail/functions/trig.hpp	Defines default versions of the trigonometric related floating point functions.

# Performance Comparison

## The Overhead in the Number Class Wrapper

Using a simple [backend class](#) that wraps any built in arithmetic type we can measure the overhead involved in wrapping a type inside the number frontend, and the effect that turning on expression templates has. The following table compares the performance between double and a double wrapped inside class number:

Type	Bessel Function Evaluation	Non-Central T Evaluation
double	<b>1.0 (0.016s)</b>	<b>1.0</b> (0.46s)
<code>number&lt;arithmetic_backend&lt;double&gt;, et_off&gt;</code>	1.2 (0.019s)	<b>1.0</b> (0.46s)
<code>number&lt;arithmetic_backend&lt;double&gt;, et_on&gt;</code>	1.2 (0.019s)	1.7 (0.79s)

As you can see whether or not there is an overhead, and how large it is depends on the actual situation, but the overhead is in any cases small. Expression templates generally add a greater overhead the more complex the expression becomes due to the logic of figuring out how to best unpack and evaluate the expression, but of course this is also the situation where you save more temporaries. For a "trivial" backend like this, saving temporaries has no benefit, but for larger types it becomes a bigger win.

The following table compares arithmetic using either `long long` or `number<arithmetic_backend<long long> >` for the [voronoi-diagram builder test](#):

Type	Relative time
<code>long long</code>	<b>1.0</b> (0.0823s)
<code>number&lt;arithmetic_backend&lt;long long&gt;, et_off&gt;</code>	1.05 (0.0875s)

This test involves mainly creating a lot of temporaries and performing a small amount of arithmetic on them, with very little difference in performance between the native and "wrapped" types.

The test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [libtommath](#)-0.42.0. The tests were run on 32-bit Windows Vista machine.

## Floating-Point Real World Tests

These tests test the total time taken to execute all of Boost.Math's test cases for these functions. In each case the best performing library gets a relative score of 1, with the total execution time given in brackets. The first three libraries listed are the various floating point types provided by this library, while for comparison, two popular C++ front-ends to [MPFR](#) ( [mpfr\\_class](#) and [mpreal](#)) are also shown.



**Table 12. Bessel Function Performance**

Library	50 Decimal Digits	100 Decimal Digits
mpfr_float	1.2 (5.78s)	1.2 (9.56s)
static_mpfr_float	1.1 (5.47s)	1.1 (9.09s)
mpf_float	<b>1.0</b> (4.82s)	<b>1.0</b> (8.07s)
cpp_dec_float	1.8 (8.54s)	2.6 (20.66s)
<a href="#">mpfr_class</a>	1.3 (6.28s)	1.2(10.06s)
<a href="#">mpreal</a>	2.0 (9.54s)	1.7 (14.08s)

**Table 13. Non-Central T Distribution Performance**

Library	50 Decimal Digits
mpfr_float	1.3 (263.27s)
static_mpfr_float	1.2 (232.88s)
mpf_float	<b>1.0</b> (195.73s)
cpp_dec_float	1.9 (366.38s)
<a href="#">mpfr_class</a>	1.5 (286.94s)
<a href="#">mpreal</a>	2.0 (388.70s)

Test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [MPFR](#)-3.0.0. The tests were run on 32-bit Windows Vista machine.

## Integer Real World Tests

The first set of [tests](#) measure the times taken to execute the multiprecision part of the Voronoi-diagram builder from Boost.Polygon. The tests mainly create a large number of temporaries "just in case" multiprecision arithmetic is required, for comparison, also included in the tests is Boost.Polygon's own partial-multiprecision integer type which was custom written for this specific task:

Integer Type	Relative Performance (Actual time in parenthesis)
polygon::detail::extended_int	1(0.138831s)
int256_t	1.19247(0.165551s)
int512_t	1.23301(0.17118s)
int1024_t	1.21463(0.168628s)
checked_int256_t	1.31711(0.182855s)
checked_int512_t	1.57413(0.218538s)
checked_int1024_t	1.36992(0.190187s)
cpp_int	1.63244(0.226632s)
mpz_int	5.42511(0.753172s)
tom_int	29.0793(4.03709s)

Note how for this use case, any dynamic allocation is a performance killer.

The next [tests](#) measure the time taken to generate 1000 128-bit random numbers and test for primality using the Miller Rabin test. This is primarily a test of modular-exponentiation since that is the rate limiting step:

Integer Type	Relative Performance (Actual time in parenthesis)
cpp_int	5.25827(0.379597s)
cpp_int (no Expression templates)	5.15675(0.372268s)
cpp_int (128-bit cache)	5.10882(0.368808s)
cpp_int (256-bit cache)	5.50623(0.397497s)
cpp_int (512-bit cache)	4.82257(0.348144s)
cpp_int (1024-bit cache)	5.00053(0.360991s)
int1024_t	4.37589(0.315897s)
checked_int1024_t	4.52396(0.326587s)
mpz_int	1(0.0721905s)
mpz_int (no Expression templates)	1.0248(0.0739806s)
tom_int	2.60673(0.188181s)
tom_int (no Expression templates)	2.64997(0.191303s)

It's interesting to note that expression templates have little effect here - perhaps because the actual expressions involved are relatively trivial in this case - so the time taken for multiplication and division tends to dominate. Also note how increasing the internal cache size used by `cpp_int` is quite effective in this case in cutting out memory allocations altogether - cutting about a third off the total

runtime. Finally the much quicker times from GMP and tommath are down to their much better modular-exponentiation algorithms (GMP's is about 5x faster). That's an issue which needs to be addressed in a future release for [cpp\\_int](#).

Test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [MPFR-3.0.0](#). The tests were run on 32-bit Windows Vista machine.

## Float Algorithm Performance

Note that these tests are carefully designed to test performance of the underlying algorithms and not memory allocation or variable copying. As usual, performance results should be taken with a healthy dose of scepticism, and real-world performance may vary widely depending upon the specifics of the program. In each table relative times are given first, with the best performer given a score of 1. Total actual times are given in brackets, measured in seconds for 500000 operations.

**Table 14. Operator +**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.0575156s)	1 (0.0740086s)	1 (0.219073s)
gmp_float	2.45065 (0.14095s)	2.01398 (0.149052s)	1.09608 (0.240122s)
mpfr_float	2.6001 (0.149546s)	2.12079 (0.156957s)	1.09078 (0.23896s)

**Table 15. Operator +(int)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1.46115 (0.0855392s)	2.60353 (0.114398s)	3.62562 (0.264905s)
gmp_float	1 (0.0585424s)	1 (0.0439398s)	1 (0.0730648s)
mpfr_float	2.40441 (0.14076s)	3.2877 (0.144461s)	2.40379 (0.175632s)

**Table 16. Operator +(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.118146s)	1 (0.144714s)	1 (0.315639s)
gmp_float	4.5555 (0.538213s)	3.83096 (0.554395s)	1.95079 (0.615745s)
mpfr_float	5.74477 (0.678719s)	4.85295 (0.702291s)	2.70354 (0.853342s)

**Table 17. Operator +=(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.101188s)	1 (0.122394s)	1 (0.251975s)
gmp_float	5.199 (0.526079s)	4.39327 (0.537712s)	2.42151 (0.610159s)
mpfr_float	6.08318 (0.615547s)	5.18525 (0.634645s)	3.1022 (0.781677s)

**Table 18. Operator -**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.0895163s)	1 (0.129248s)	1.5088 (0.374512s)
gmp_float	1.72566 (0.154474s)	1.22567 (0.158415s)	1 (0.248219s)
mpfr_float	1.83764 (0.164499s)	1.34284 (0.173559s)	1.00226 (0.248781s)

**Table 19. Operator -(int)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.105285s)	1 (0.142741s)	1 (0.278718s)
gmp_float	2.34437 (0.246828s)	1.28814 (0.183871s)	1.00731 (0.280754s)
mpfr_float	2.8032 (0.295136s)	2.09178 (0.298582s)	1.25213 (0.34899s)

**Table 20. Operator -(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.13719s)	1 (0.184428s)	1 (0.344212s)
gmp_float	4.0804 (0.559791s)	3.06776 (0.565781s)	2.07736 (0.715053s)
mpfr_float	5.10114 (0.699828s)	3.88684 (0.716843s)	2.50074 (0.860784s)

**Table 21. Operator ==(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.100984s)	1 (0.123148s)	1 (0.246181s)
gmp_float	5.68353 (0.573944s)	4.68636 (0.577116s)	2.6958 (0.663655s)
mpfr_float	6.19738 (0.625834s)	5.18544 (0.638577s)	3.18738 (0.784673s)

**Table 22. Operator \***

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1.03667 (0.284251s)	1.30576 (0.536527s)	1.44686 (4.81057s)
gmp_float	1 (0.274196s)	1 (0.410891s)	1 (3.32484s)
mpfr_float	1.24537 (0.341477s)	1.15785 (0.475749s)	1.1796 (3.92199s)

**Table 23. Operator \*(int)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	3.97453 (0.240262s)	9.91222 (0.463473s)	50.7926 (4.36527s)
gmp_float	<b>1</b> (0.0604505s)	<b>1</b> (0.0467577s)	<b>1</b> (0.0859431s)
mpfr_float	2.56974 (0.155342s)	3.56312 (0.166603s)	3.22964 (0.277565s)

**Table 24. Operator \*(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	<b>1</b> (0.331877s)	1.01058 (0.586122s)	6.688 (4.7931s)
gmp_float	1.72433 (0.572266s)	<b>1</b> (0.579987s)	<b>1</b> (0.716672s)
mpfr_float	2.5553 (0.848047s)	1.74987 (1.0149s)	1.80403 (1.2929s)

**Table 25. Operator \*=(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	<b>1</b> (0.321397s)	1.00772 (0.574887s)	6.65946 (4.7468s)
gmp_float	1.77419 (0.570218s)	<b>1</b> (0.570482s)	<b>1</b> (0.712791s)
mpfr_float	2.62172 (0.842611s)	1.77691 (1.01369s)	1.77511 (1.26528s)

**Table 26. Operator /**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	2.96096 (4.00777s)	4.53244 (7.86435s)	6.11936 (51.5509s)
gmp_float	<b>1</b> (1.35354s)	<b>1</b> (1.73512s)	<b>1</b> (8.42422s)
mpfr_float	1.30002 (1.75963s)	1.39045 (2.41261s)	1.66762 (14.0484s)

**Table 27. Operator /(int)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	8.60726 (1.8181s)	15.4122 (3.67479s)	34.5119 (24.729s)
gmp_float	1.24394 (0.262756s)	<b>1</b> (0.238433s)	<b>1</b> (0.716536s)
mpfr_float	<b>1</b> (0.211229s)	1.12178 (0.26747s)	1.02237 (0.732562s)

**Table 28. Operator /(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	2.10976 (1.97569s)	3.73601 (3.9133s)	11.3085 (25.4533s)
gmp_float	1 (0.936452s)	1 (1.04746s)	1 (2.25081s)
mpfr_float	1.3423 (1.257s)	1.51575 (1.58768s)	3.31513 (7.46175s)

**Table 29. Operator /=(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	2.17401 (1.96883s)	3.79591 (3.8965s)	11.2328 (25.2606s)
gmp_float	1 (0.905621s)	1 (1.0265s)	1 (2.24882s)
mpfr_float	1.37953 (1.24933s)	1.53073 (1.57129s)	3.30546 (7.43339s)

**Table 30. Operator construct**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.00929804s)	1 (0.0268321s)	1 (0.0310685s)
gmp_float	30.8781 (0.287106s)	7.59969 (0.203916s)	6.51873 (0.202527s)
mpfr_float	23.5296 (0.218779s)	8.11058 (0.217624s)	7.16325 (0.222552s)

**Table 31. Operator construct(unsigned)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.0603971s)	1 (0.0735485s)	1 (0.116464s)
gmp_float	3.91573 (0.236498s)	2.88171 (0.211945s)	1.81075 (0.210887s)
mpfr_float	4.90052 (0.295977s)	4.01118 (0.295017s)	2.62005 (0.305141s)

**Table 32. Operator construct(unsigned long long)**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	1 (0.0610288s)	1 (0.0759005s)	1 (0.118511s)
gmp_float	8.26247 (0.504249s)	6.69042 (0.507806s)	4.32819 (0.51294s)
mpfr_float	10.1593 (0.620013s)	8.45884 (0.64203s)	5.51472 (0.653557s)

**Table 33. Operator str**

Backend	50 Bits	100 Bits	500 Bits
cpp_dec_float	2.95848 (0.0223061s)	3.33461 (0.033471s)	3.0159 (0.132732s)
gmp_float	<b>1</b> (0.00753971s)	<b>1</b> (0.0100374s)	<b>1</b> (0.0440106s)
mpfr_float	1.25424 (0.00945658s)	1.24943 (0.012541s)	1.09428 (0.0481601s)

Test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [MPFR-3.0.0](#). The tests were run on 32-bit Windows Vista machine.

## Integer Algorithm Performance

Note that these tests are carefully designed to test performance of the underlying algorithms and not memory allocation or variable copying. As usual, performance results should be taken with a healthy dose of scepticism, and real-world performance may vary widely depending upon the specifics of the program. In each table relative times are given first, with the best performer given a score of 1. Total actual times are given in brackets, measured in seconds for 500000 operations.

**Table 34. Operator +**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.23704 (0.0274266s)	1.09358 (0.0383278s)	1.26645 (0.0558828s)	1.32188 (0.0916899s)
cpp_int(fixed)	1.62044 (0.0359271s)	1.5277 (0.053543s)	1.73059 (0.076363s)	1.71537 (0.118983s)
gmp_int	1.87515 (0.0415741s)	1.21699 (0.042653s)	1.15599 (0.0510088s)	<b>1</b> (0.0693631s)
tommath_int	<b>1</b> (0.0221711s)	<b>1</b> (0.035048s)	<b>1</b> (0.0441255s)	1.04441 (0.0724435s)

**Table 35. Operator +(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0155377s)	<b>1</b> (0.0209523s)	<b>1</b> (0.0306377s)	<b>1</b> (0.043125s)
cpp_int(fixed)	1.31904 (0.0204948s)	1.76211 (0.0369203s)	1.52941 (0.0468577s)	1.60412 (0.0691778s)
gmp_int	1.96204 (0.0304855s)	2.02569 (0.0424428s)	2.11505 (0.0648002s)	2.65993 (0.114709s)
tommath_int	14.0654 (0.218543s)	10.8239 (0.226786s)	7.76691 (0.23796s)	6.10039 (0.263079s)

**Table 36. Operator +(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.026624s)	<b>1</b> (0.0291407s)	<b>1</b> (0.0373209s)	<b>1</b> (0.0464919s)
cpp_int(fixed)	1.31378 (0.034978s)	1.54897 (0.045138s)	1.53649 (0.0573431s)	1.27833 (0.0594319s)
gmp_int	25.5775 (0.680974s)	24.0117 (0.699717s)	19.5633 (0.730121s)	16.8939 (0.785432s)
tommath_int	19.4694 (0.518354s)	18.4246 (0.536907s)	14.7715 (0.551288s)	12.3637 (0.574812s)

**Table 37. Operator +=(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.18405 (0.0196905s)	1.22304 (0.0206476s)	1.25861 (0.0217397s)	1.29525 (0.0220829s)
cpp_int(fixed)	<b>1</b> (0.0166298s)	<b>1</b> (0.0168822s)	<b>1</b> (0.0172728s)	<b>1</b> (0.0170492s)
gmp_int	39.9082 (0.663668s)	39.4584 (0.666147s)	38.5504 (0.665873s)	39.2231 (0.668722s)
tommath_int	30.6219 (0.509238s)	30.4135 (0.513447s)	30.9077 (0.533863s)	32.3086 (0.550835s)

**Table 38. Operator -**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.06986 (0.0296064s)	<b>1</b> (0.0381508s)	1.05932 (0.053186s)	1.1766 (0.0844721s)
cpp_int(fixed)	1.3304 (0.0368163s)	1.44506 (0.0551303s)	1.4431 (0.0724545s)	1.57255 (0.112898s)
gmp_int	1.48072 (0.0409761s)	1.19003 (0.0454007s)	1.0794 (0.0541942s)	<b>1</b> (0.0717934s)
tommath_int	<b>1</b> (0.0276731s)	1.10891 (0.0423057s)	<b>1</b> (0.0502076s)	1.08479 (0.0778811s)

**Table 39. Operator -(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0147372s)	<b>1</b> (0.0170001s)	<b>1</b> (0.0232882s)	<b>1</b> (0.0310734s)
cpp_int(fixed)	1.4267 (0.0210256s)	1.98887 (0.0338109s)	1.83788 (0.0428009s)	1.81269 (0.0563264s)
gmp_int	2.07504 (0.0305803s)	2.40928 (0.0409579s)	2.58711 (0.0602493s)	3.26438 (0.101435s)
tommath_int	13.5424 (0.199577s)	12.1793 (0.207048s)	9.28855 (0.216314s)	7.49327 (0.232842s)

**Table 40. Operator -(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0277377s)	<b>1</b> (0.0296807s)	<b>1</b> (0.0372392s)	<b>1</b> (0.0455855s)
cpp_int(fixed)	1.19867 (0.0332484s)	1.48639 (0.0441169s)	1.43253 (0.0533464s)	1.27697 (0.0582111s)
gmp_int	24.1794 (0.670683s)	22.9073 (0.679904s)	18.8758 (0.702922s)	16.5837 (0.755975s)
tommath_int	18.149 (0.503413s)	17.4116 (0.516787s)	14.0411 (0.52288s)	11.8237 (0.538987s)



**Table 41. Operator  $\text{--}=(\text{unsigned long long})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.26896 (0.0203467s)	1.25722 (0.0206147s)	1.36108 (0.0225485s)	1.18351 (0.0226161s)
cpp_int(fixed)	<b>1</b> (0.0160342s)	<b>1</b> (0.0163971s)	<b>1</b> (0.0165667s)	<b>1</b> (0.0191094s)
gmp_int	41.1339 (0.659547s)	40.3982 (0.662411s)	39.925 (0.661425s)	34.636 (0.661874s)
tommath_int	31.1543 (0.499533s)	31.0303 (0.508806s)	30.7699 (0.509756s)	27.7054 (0.529434s)

**Table 42. Operator  $*$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.11839 (0.0757577s)	1.61061 (0.207951s)	1.4501 (0.696912s)	1.72796 (2.64108s)
cpp_int(fixed)	1.01115 (0.0684934s)	1.28687 (0.166152s)	<b>1</b> (0.480595s)	<b>1</b> (1.52844s)
gmp_int	<b>1</b> (0.0677384s)	<b>1</b> (0.129113s)	1.09011 (0.523902s)	1.03374 (1.58s)
tommath_int	1.6322 (0.110562s)	2.71751 (0.350866s)	2.05222 (0.986288s)	2.0644 (3.15531s)

**Table 43. Operator  $*(\text{int})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.01611 (0.0229536s)	1.12175 (0.0298152s)	1.16413 (0.0416439s)	1.31747 (0.0666043s)
cpp_int(fixed)	1.30215 (0.0294152s)	1.669 (0.0443606s)	1.72395 (0.0616701s)	1.88315 (0.095202s)
gmp_int	<b>1</b> (0.0225897s)	<b>1</b> (0.0265791s)	<b>1</b> (0.0357725s)	<b>1</b> (0.0505547s)
tommath_int	10.8281 (0.244603s)	10.1516 (0.26982s)	8.76424 (0.313519s)	8.04364 (0.406644s)

**Table 44. Operator  $*(\text{unsigned long long})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0570721s)	<b>1</b> (0.0856141s)	<b>1</b> (0.143279s)	<b>1</b> (0.252785s)
cpp_int(fixed)	1.10857 (0.0632686s)	1.2951 (0.110878s)	1.20827 (0.173121s)	1.18463 (0.299456s)
gmp_int	12.0605 (0.68832s)	8.13434 (0.696415s)	5.21762 (0.747577s)	3.11601 (0.787681s)
tommath_int	10.0524 (0.57371s)	7.33116 (0.627651s)	4.85202 (0.695193s)	3.35808 (0.848871s)

**Table 45. Operator  $\ast=(\text{unsigned long long})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	111.27 (7.43118s)	67.7078 (7.34138s)	43.3851 (7.4075s)	25.3089 (7.55455s)
cpp_int(fixed)	<b>1</b> (0.0667848s)	<b>1</b> (0.108427s)	<b>1</b> (0.170738s)	<b>1</b> (0.298493s)
gmp_int	46.3718 (3.09693s)	28.4639 (3.08626s)	18.1719 (3.10264s)	10.5223 (3.14083s)
tommath_int	276.674 (18.4776s)	169.146 (18.34s)	108.491 (18.5236s)	63.3261 (18.9024s)

**Table 46. Operator  $/$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	2.68035 (0.595251s)	2.04702 (0.707471s)	1.62314 (0.921536s)	1.43112 (1.38811s)
cpp_int(fixed)	<b>1</b> (0.222079s)	<b>1</b> (0.34561s)	<b>1</b> (0.567748s)	<b>1</b> (0.969945s)
gmp_int	3.79283 (0.842308s)	2.73668 (0.945824s)	1.86649 (1.05969s)	1.32141 (1.2817s)
tommath_int	13.2531 (2.94324s)	11.2054 (3.87271s)	9.83293 (5.58262s)	13.0164 (12.6252s)

**Table 47. Operator  $/(\text{int})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	4.06026 (0.225473s)	3.45732 (0.340049s)	3.00195 (0.547957s)	2.80587 (0.978029s)
cpp_int(fixed)	2.43766 (0.135367s)	2.56264 (0.252052s)	2.44011 (0.445402s)	2.38009 (0.829617s)
gmp_int	<b>1</b> (0.0555316s)	<b>1</b> (0.0983563s)	<b>1</b> (0.182534s)	<b>1</b> (0.348566s)
tommath_int	35.9988 (1.99907s)	27.1024 (2.66569s)	21.8333 (3.98531s)	25.8066 (8.99528s)

**Table 48. Operator  $/(\text{unsigned long long})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.50505 (0.705756s)	1.39347 (1.58556s)	2.63348 (3.57438s)	4.75451 (8.52733s)
cpp_int(fixed)	<b>1</b> (0.468925s)	1.12378 (1.27869s)	2.29966 (3.12128s)	4.4844 (8.04288s)
gmp_int	2.17234 (1.01866s)	<b>1</b> (1.13785s)	<b>1</b> (1.35728s)	<b>1</b> (1.79352s)
tommath_int	4.74612 (2.22557s)	2.70088 (3.07319s)	3.65634 (4.96268s)	6.79408 (12.1853s)

**Table 49. Operator /= (unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.76281 (0.0574966s)	1.76471 (0.0604224s)	1.56085 (0.0716403s)	1.31422 (0.124043s)
cpp_int(fixed)	<b>1</b> (0.0326164s)	<b>1</b> (0.0342393s)	<b>1</b> (0.0458981s)	<b>1</b> (0.0943852s)
gmp_int	20.2862 (0.661664s)	19.4043 (0.664389s)	14.4881 (0.664976s)	7.14238 (0.674135s)
tommath_int	32.9555 (1.07489s)	30.1525 (1.0324s)	22.8324 (1.04796s)	11.7456 (1.10861s)

**Table 50. Operator %**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.8501 (0.364131s)	1.46527 (0.476653s)	1.27509 (0.689738s)	1.20064 (1.11769s)
cpp_int(fixed)	<b>1</b> (0.196817s)	<b>1</b> (0.325301s)	<b>1</b> (0.540932s)	<b>1</b> (0.930916s)
gmp_int	3.2533 (0.640305s)	2.15441 (0.700832s)	1.47898 (0.800029s)	1.07439 (1.00016s)
tommath_int	15.3501 (3.02116s)	12.1106 (3.9396s)	11.0689 (5.98752s)	13.5535 (12.6172s)

**Table 51. Operator %(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.82761 (0.104331s)	2.01496 (0.202512s)	2.10004 (0.389523s)	2.17252 (0.768097s)
cpp_int(fixed)	1.78851 (0.102099s)	1.96844 (0.197838s)	2.02956 (0.376451s)	2.07257 (0.73276s)
gmp_int	<b>1</b> (0.057086s)	<b>1</b> (0.100505s)	<b>1</b> (0.185483s)	<b>1</b> (0.353552s)
tommath_int	36.3018 (2.07233s)	26.3075 (2.64402s)	21.9525 (4.07183s)	25.6759 (9.07775s)

**Table 52. Operator construct**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.40211 (0.0026854s)	<b>1</b> (0.00278639s)	<b>1</b> (0.00322813s)	<b>1</b> (0.0027185s)
cpp_int(fixed)	<b>1</b> (0.00191526s)	1.40721 (0.00392103s)	1.90346 (0.00614463s)	2.14621 (0.00583447s)
gmp_int	98.705 (0.189046s)	68.9726 (0.192184s)	58.8994 (0.190135s)	70.0525 (0.190438s)
tommath_int	105.602 (0.202255s)	74.1994 (0.206748s)	63.6455 (0.205456s)	76.8935 (0.209035s)

**Table 53. Operator construct(unsigned)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.73436 (0.00348927s)	1 (0.00263476s)	1 (0.0027009s)	1 (0.00318651s)
cpp_int(fixed)	1 (0.00201185s)	1.36851 (0.0036057s)	2.07362 (0.00560064s)	1.66856 (0.00531688s)
gmp_int	97.2414 (0.195635s)	76.3759 (0.201232s)	72.7396 (0.196462s)	63.8129 (0.20334s)
tommath_int	210.112 (0.422713s)	162.652 (0.42855s)	158.33 (0.427634s)	134.626 (0.428987s)

**Table 54. Operator construct(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	2.34403 (0.00739542s)	1.66376 (0.00713834s)	1.22989 (0.0074969s)	1.23708 (0.00711417s)
cpp_int(fixed)	1 (0.00315501s)	1 (0.00429049s)	1 (0.00609561s)	1 (0.0057508s)
gmp_int	222.866 (0.703144s)	164.331 (0.705059s)	115.363 (0.70321s)	122.347 (0.703596s)
tommath_int	218.681 (0.689941s)	163.796 (0.702765s)	114.57 (0.698376s)	122.422 (0.704027s)

**Table 55. Operator gcd**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.16358 (2.74442s)	1.39847 (8.11559s)	1.64677 (22.2518s)	1.95096 (64.4961s)
cpp_int(fixed)	1 (2.35859s)	1.30986 (7.60133s)	1.67681 (22.6577s)	2.0895 (69.0758s)
gmp_int	1.03392 (2.4386s)	1 (5.80319s)	1 (13.5124s)	1 (33.0586s)
tommath_int	5.25978 (12.4057s)	4.4619 (25.8932s)	4.15577 (56.1542s)	3.91192 (129.323s)

**Table 56. Operator powm**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	2.50722 (2.91621s)	3.5561 (13.406s)	4.37066 (73.483s)	4.88831 (473.91s)
cpp_int(fixed)	1.93385 (2.24931s)	3.18107 (11.9922s)	4.20753 (70.7403s)	4.8158 (466.88s)
gmp_int	1 (1.16313s)	1 (3.76986s)	1 (16.8128s)	1 (96.9476s)
tommath_int	1.44081 (1.67584s)	1.8794 (7.08507s)	2.19115 (36.8394s)	2.17186 (210.557s)

**Table 57. Operator str**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.17175 (0.00160006s)	1.41999 (0.00329476s)	1.40856 (0.00813784s)	1.52964 (0.0229767s)
cpp_int(fixed)	<b>1</b> (0.00136554s)	<b>1</b> (0.00232027s)	<b>1</b> (0.00577741s)	1.14754 (0.0172372s)
gmp_int	1.50501 (0.00205515s)	1.52968 (0.00354926s)	1.01989 (0.0058923s)	<b>1</b> (0.015021s)
tommath_int	12.2161 (0.0166816s)	16.9577 (0.0393463s)	18.7474 (0.108311s)	22.7368 (0.341528s)

**Table 58. Operator |**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0301617s)	<b>1</b> (0.0423404s)	<b>1</b> (0.0522358s)	<b>1</b> (0.0813156s)
cpp_int(fixed)	1.0638 (0.0320861s)	1.22566 (0.0518951s)	1.28515 (0.0671305s)	1.16118 (0.094422s)
gmp_int	1.76553 (0.0532514s)	1.51489 (0.0641408s)	1.70708 (0.0891706s)	1.77346 (0.14421s)
tommath_int	4.37637 (0.131999s)	3.46212 (0.146587s)	2.91875 (0.152463s)	4.19621 (0.341217s)

**Table 59. Operator |(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0289129s)	<b>1</b> (0.0351119s)	<b>1</b> (0.0406779s)	<b>1</b> (0.0525891s)
cpp_int(fixed)	1.06091 (0.030674s)	1.25979 (0.0442336s)	1.36194 (0.0554009s)	1.37438 (0.0722772s)
gmp_int	4.92854 (0.142498s)	4.34687 (0.152627s)	3.71442 (0.151095s)	2.981 (0.156768s)
tommath_int	10.9847 (0.317598s)	9.37065 (0.329021s)	8.53651 (0.347248s)	11.2155 (0.589813s)

**Table 60. Operator ^**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	<b>1</b> (0.0305149s)	<b>1</b> (0.04217s)	<b>1</b> (0.0525977s)	<b>1</b> (0.0816632s)
cpp_int(fixed)	1.01544 (0.0309861s)	1.24872 (0.0526585s)	1.26661 (0.066621s)	1.15965 (0.0947007s)
gmp_int	1.64675 (0.0502505s)	1.47181 (0.0620663s)	1.66038 (0.0873322s)	1.67895 (0.137108s)
tommath_int	4.30668 (0.131418s)	3.45859 (0.145849s)	2.91462 (0.153303s)	4.15538 (0.339342s)

**Table 61. Operator  $^{\wedge}(\text{int})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.01566 (0.0296088s)	1 (0.0356634s)	1 (0.0401898s)	1 (0.0514097s)
cpp_int(fixed)	1 (0.0291524s)	1.2393 (0.0441976s)	1.38556 (0.0556856s)	1.38899 (0.0714075s)
gmp_int	4.68027 (0.136441s)	4.15243 (0.14809s)	3.74237 (0.150405s)	3.0483 (0.156712s)
tommath_int	10.919 (0.318314s)	9.16311 (0.326788s)	8.62554 (0.346659s)	11.6212 (0.597442s)

**Table 62. Operator  $\&$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.0346 (0.0303431s)	1 (0.0427309s)	1 (0.0535587s)	1.06945 (0.0828084s)
cpp_int(fixed)	1 (0.0293284s)	1.10435 (0.04719s)	1.05262 (0.0563769s)	1 (0.0774309s)
gmp_int	1.86057 (0.0545675s)	1.58432 (0.0676995s)	1.69164 (0.0906018s)	1.86625 (0.144505s)
tommath_int	4.4157 (0.129506s)	3.60396 (0.154s)	2.95985 (0.158525s)	4.4032 (0.340944s)

**Table 63. Operator  $\&(\text{int})$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1.05874 (0.038946s)	1 (0.0483903s)	1 (0.063842s)	1 (0.100361s)
cpp_int(fixed)	1 (0.0367853s)	1.05827 (0.0512099s)	1.09114 (0.0696605s)	1.09432 (0.109826s)
gmp_int	3.92298 (0.144308s)	2.99447 (0.144903s)	2.228 (0.14224s)	1.42296 (0.142809s)
tommath_int	8.79208 (0.323419s)	7.02288 (0.339839s)	5.65271 (0.36088s)	6.27104 (0.629365s)

**Table 64. Operator  $\ll$** 

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1 (0.0248801s)	1.23196 (0.04s)	1 (0.0424149s)	1 (0.060157s)
cpp_int(fixed)	1.08931 (0.027102s)	1.40572 (0.0456418s)	1.3475 (0.0571542s)	1.24573 (0.0749397s)
gmp_int	1.05561 (0.0262636s)	1 (0.0324686s)	1.09914 (0.0466199s)	1.16315 (0.0699719s)
tommath_int	1.60497 (0.0399319s)	2.13048 (0.0691737s)	2.31219 (0.0980712s)	2.74695 (0.165248s)

**Table 65. Operator >>**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_int	1 (0.0213349s)	1.02127 (0.0295019s)	1 (0.0327116s)	1.13168 (0.0433804s)
cpp_int(fixed)	1.13514 (0.0242181s)	1.16938 (0.0337803s)	1.46999 (0.0480859s)	1.60077 (0.061362s)
gmp_int	1.26614 (0.0270129s)	1 (0.0288873s)	1.42219 (0.0465221s)	1 (0.0383329s)
tommath_int	12.0066 (0.25616s)	10.2837 (0.297067s)	9.99696 (0.327017s)	16.0943 (0.616942s)

Test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [MPFR-3.0.0](#). The tests were run on 32-bit Windows Vista machine.

Linux x86\_64 results are broadly similar, except that libtommath performs much better there.

## Rational Type Performance

Note that these tests are carefully designed to test performance of the underlying algorithms and not memory allocation or variable copying. As usual, performance results should be taken with a healthy dose of scepticism, and real-world performance may vary widely depending upon the specifics of the program. In each table relative times are given first, with the best performer given a score of 1. Total actual times are given in brackets, measured in seconds for 500000 operations.

**Table 66. Operator +**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	5.89417 (18.4116s)	6.87256 (47.4698s)	6.65008 (107.715s)	6.53801 (256.244s)
mpq_rational	1 (3.1237s)	1 (6.90715s)	1 (16.1975s)	1 (39.1929s)

**Table 67. Operator +(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	3.62367 (2.46488s)	4.18291 (2.94603s)	4.726 (3.74866s)	6.1388 (5.56817s)
mpq_rational	1 (0.680215s)	1 (0.704303s)	1 (0.7932s)	1 (0.907046s)

**Table 68. Operator +(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.1527 (2.6378s)	1.31751 (3.09863s)	1.58996 (4.00714s)	2.15642 (5.75702s)
mpq_rational	1 (2.28837s)	1 (2.35189s)	1 (2.52028s)	1 (2.66971s)

**Table 69. Operator +=(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.18436 (2.7059s)	1.32279 (3.11099s)	1.61398 (4.05389s)	2.20048 (5.84623s)
mpq_rational	<b>1</b> (2.2847s)	<b>1</b> (2.35183s)	<b>1</b> (2.51174s)	<b>1</b> (2.6568s)

**Table 70. Operator -**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	5.81893 (18.3457s)	6.82209 (47.1928s)	6.64143 (107.498s)	6.51362 (255.137s)
mpq_rational	<b>1</b> (3.15277s)	<b>1</b> (6.91765s)	<b>1</b> (16.1859s)	<b>1</b> (39.1698s)

**Table 71. Operator -(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	3.72441 (2.48756s)	4.27663 (2.98713s)	4.62109 (3.72114s)	6.17605 (5.56503s)
mpq_rational	<b>1</b> (0.667908s)	<b>1</b> (0.698479s)	<b>1</b> (0.805252s)	<b>1</b> (0.901066s)

**Table 72. Operator -(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.15627 (2.63239s)	1.32096 (3.12092s)	1.61044 (4.00106s)	2.19378 (5.7644s)
mpq_rational	<b>1</b> (2.27663s)	<b>1</b> (2.36262s)	<b>1</b> (2.48445s)	<b>1</b> (2.62761s)

**Table 73. Operator -=(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.1984 (2.73444s)	1.34141 (3.15698s)	1.64159 (4.06997s)	2.23017 (5.88108s)
mpq_rational	<b>1</b> (2.28174s)	<b>1</b> (2.35348s)	<b>1</b> (2.47929s)	<b>1</b> (2.63706s)

**Table 74. Operator \***

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	5.4306 (32.5882s)	6.91805 (89.9436s)	6.94556 (207.307s)	6.88704 (492.151s)
mpq_rational	<b>1</b> (6.00084s)	<b>1</b> (13.0013s)	<b>1</b> (29.8475s)	<b>1</b> (71.4604s)



**Table 75. Operator \*(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	2.12892 (2.51376s)	2.47245 (3.07841s)	2.86832 (3.93619s)	3.94086 (6.02565s)
mpq_rational	<b>1</b> (1.18077s)	<b>1</b> (1.24508s)	<b>1</b> (1.3723s)	<b>1</b> (1.52902s)

**Table 76. Operator \*(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.32254 (5.43565s)	1.56078 (6.73163s)	1.97701 (9.32522s)	2.85404 (15.1573s)
mpq_rational	<b>1</b> (4.11002s)	<b>1</b> (4.313s)	<b>1</b> (4.71682s)	<b>1</b> (5.31082s)

**Table 77. Operator \*=(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	6.29806 (58.1188s)	6.30556 (59.5076s)	6.3385 (62.1007s)	6.55345 (67.6905s)
mpq_rational	<b>1</b> (9.22804s)	<b>1</b> (9.43733s)	<b>1</b> (9.79739s)	<b>1</b> (10.329s)

**Table 78. Operator /**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	4.4269 (66.8031s)	6.40103 (173.527s)	6.32347 (348.193s)	6.61148 (824.063s)
mpq_rational	<b>1</b> (15.0903s)	<b>1</b> (27.1093s)	<b>1</b> (55.0637s)	<b>1</b> (124.641s)

**Table 79. Operator /(int)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.78772 (2.50984s)	2.10623 (3.10606s)	2.46986 (3.99358s)	3.37428 (5.96678s)
mpq_rational	<b>1</b> (1.40393s)	<b>1</b> (1.4747s)	<b>1</b> (1.61693s)	<b>1</b> (1.76831s)

**Table 80. Operator /(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1.29695 (5.45454s)	1.55248 (6.85353s)	1.93237 (9.28765s)	2.75211 (14.8541s)
mpq_rational	<b>1</b> (4.20568s)	<b>1</b> (4.41458s)	<b>1</b> (4.80635s)	<b>1</b> (5.39734s)

**Table 81. Operator /=(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	6.19401 (58.4278s)	6.20135 (59.643s)	6.21327 (62.0338s)	6.40576 (67.6778s)
mpq_rational	1 (9.43295s)	1 (9.61774s)	1 (9.98407s)	1 (10.5652s)

**Table 82. Operator construct**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1 (0.00978288s)	1 (0.0100574s)	1 (0.0101393s)	1 (0.0101847s)
mpq_rational	39.1516 (0.383015s)	38.3523 (0.385725s)	37.5812 (0.381048s)	37.6007 (0.382953s)

**Table 83. Operator construct(unsigned)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1 (0.0548151s)	1 (0.0557542s)	1 (0.055825s)	1 (0.0552808s)
mpq_rational	7.21073 (0.395257s)	7.1016 (0.395944s)	7.02046 (0.391917s)	7.16881 (0.396297s)

**Table 84. Operator construct(unsigned long long)**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	1 (0.0605156s)	1 (0.0616657s)	1 (0.0592056s)	1 (0.0603081s)
mpq_rational	35.1604 (2.12775s)	34.7575 (2.14335s)	35.7232 (2.11502s)	35.0437 (2.11342s)

**Table 85. Operator str**

Backend	128 Bits	256 Bits	512 Bits	1024 Bits
cpp_rational	5.48898 (0.0208949s)	8.49668 (0.0546688s)	10.107 (0.121897s)	10.5339 (0.310584s)
mpq_rational	1 (0.0038067s)	1 (0.00643413s)	1 (0.0120606s)	1 (0.0294843s)

Test code was compiled with Microsoft Visual Studio 2010 with all optimisations turned on (/Ox), and used MPIR-2.3.0 and [MPFR-3.0.0](#). The tests were run on 32-bit Windows Vista machine.

# Roadmap

## History

### Multiprecision-2.2.2 (Boost-1.56)

- Change floating point to rational conversions to be implicit, see [10082](#).
- Fix definition of `checked_cpp_rational` typedef.

### Multiprecision-2.2.1

- Fix bug in assignment from string in `cpp_int`, see [9936](#).

### Multiprecision-2.2.0

- Moved to Boost.Multiprecision specific version number - we have one breaking change in Boost-1.54 which makes this major version 2, plus two releases with new features since then.
- Added new `cpp_bin_float` backend for binary floating point.
- Added MSVC-specific `#include` for compiler intrinsics, see [9336](#).
- Fixed various typos in docs, see [9432](#).
- Fixed `gmp_rational` to allow move-copy from an already copied-from object, see [9497](#).
- Added list of values for `numeric_limits`.

### Boost-1.55

- Added support for Boost.Serialization.
- Suppressed some GCC warnings. See [8872](#).
- Fixed bug in `pow` for large integer arguments. See [8809](#).
- Fixed bug in `pow` for calculation of  $0^N$ . See [8798](#).
- Fixed bug in fixed precision `cpp_int` IO code that causes conversion to string to fail when the bit count is very small (less than `CHAR_BIT`). See [8745](#).
- Fixed bug in `cpp_int` that causes left shift to fail when a fixed precision type would overflow. See [8741](#).
- Fixed some cosmetic warnings from `cpp_int`. See [8748](#).
- Fixed calls to functions which are required to be macros in C99. See [8732](#).
- Fixed bug that causes construction from `INT_MIN`, `LONG_MIN` etc to fail in `cpp_int`. See [8711](#).

### 1.54

- **Breaking change** renamed `rational_adapter` to `rational_adaptor`.
- Add support for [MPFI](#).
- Add `logged_adaptor`.
- Add support for 128-bit floats via GCC's `float128` or Intel's `_Quad` data types.

- Add support for user-defined literals in `cpp_int`, improve `constexpr` support.
- Fixed bug in integer division of `cpp_int` that results in incorrect sign of `cpp_int` when both arguments are small enough to fit in a `double_limb_type`. See [8126](#).
- Fixed bug in subtraction of a single limb in `cpp_int` that results in incorrect value when the result should have a 0 in the last limb: [8133](#).
- Fixed bug in `cpp_int` where division of 0 by something doesn't get zero in the result: [8160](#).
- Fixed bug in some transcendental functions that caused incorrect return values when variables are reused, for example with `a = pow(a, b)`. See [8326](#).
- Fixed some assignment operations in the `mpfr` and `gmp` backends to be safe if the target has been moved from: [8667](#).
- Fixed bug in `cpp_int` that gives incorrect answer for `0%N` for large `N`: [8670](#).
- Fixed `set_precision` in `mpfr` backend so it doesn't trample over an existing value: [8692](#).

## 1.53

- First Release.
- Fix bug in `cpp_int` division.
- Fix issue [#7806](#).

## Post review changes

- Non-expression template operators further optimised with rvalue reference support.
- Many functions made `constexpr`.
- Differentiate between explicit and implicit conversions in the number constructor.
- Removed "mp\_" prefix from types.
- Allowed mixed precision arithmetic.
- Changed `ExpressionTemplates` parameter to class `number` to use enumerated values rather than `true/false`.
- Changed `ExpressionTemplate` parameter default value to use a traits class so that the default value depends on the backend used.
- Added support for fused-multiply-add/subtract with GMP support.
- Tweaked expression template unpacking to use fewer temporaries when the LHS also appears in the RHS.
- Refactored `cpp_int_backend` based on review comments with new template parameter structure.
- Added additional template parameter to `mpfr_float_backend` to allow stack-based allocation.
- Added section on mixed precision arithmetic, and added support for operations yielding a higher precision result than either of the arguments.
- Added overloads of integer-specific functions for built in integer types.

## Pre-review history

- 2011-2012, John Maddock adds an expression template enabled front end to Christopher's code, and adds support for other backends.

- 2011, Christopher Kormanyos publishes the decimal floating point code under the Boost Software Licence. The code is published as: "[Algorithm 910: A Portable C++ Multiple-Precision System for Special-Function Calculations](#)", in ACM TOMS, {VOL 37, ISSUE 4, (February 2011)} (C) ACM, 2011.
- 2002-2011, Christopher Kormanyos develops the all C++ decimal arithmetic floating point code.

## TODO

More a list of what *could* be done, rather than what *should* be done (which may be a much smaller list!).

- Add back-end support for libdecNumber.
- Add an adaptor back-end for complex number types.
- Add a back-end for MPFR interval arithmetic.
- Add better multiplication routines (Karatsuba, FFT etc) to `cpp_int_backend`.
- Add assembly level routines to `cpp_int_backend`.
- Add an all C++ Boost licensed binary floating point type.
- Can ring types (exact floating point types) be supported? The answer should be yes, but someone needs to write it, the hard part is IO and binary-decimal conversion.
- Should there be a choice of rounding mode (probably MPFR specific)?
- We can reuse temporaries in multiple subtrees (temporary caching).
- `cpp_dec_float` should round to nearest.
- A 2's complement fixed precision int that uses exactly N bits and no more.

Things requested in review:

- The performances of `mp_number<a_trivial_adaptor<float>, false>` respect to float and `mp_number<a_trivial_adaptor<int>, false>` and int should be given to show the cost of using the generic interface (Mostly done, just need to update docs to the latest results).
- Should we provide min/max overloads for expression templates? (Not done - we can't overload functions declared in the std namespace :-()).
- The rounding applied when converting must be documented (Done).
- Document why we don't abstract out addition/multiplication algorithms etc. (done - FAQ)
- Document why we don't use proto (compile times) (Done).
- We can reuse temporaries in multiple subtrees (temporary caching) Moved to TODO list.
- Emphasise in the docs that ET's may reorder operations (done 2012/10/31).
- Document what happens to small fixed precision `cpp_int`'s (done 2012/10/31).
- The use of bool in template parameters could be improved by the use of an enum class which will be more explicit. E.g `enum class expression_template {disabled, enabled}; enum class sign {unsigned, signed};` (Partly done 2012/09/15, done 2012/10/31).
- Each back-end should document the requirements it satisfies (not currently scheduled for inclusion: it's deliberately an implementation detail, and "optional" requirements are optimisations which can't be detected by the user). Not done: this is an implementation detail, the exact list of requirements satisfied is purely an optimization, not something the user can detect.

- A backend for an overflow aware integers (done 2012/10/31).
- IIUC `convert_to` is used to emulate in c++98 compilers C++11 explicit conversions. Could the explicit conversion operator be added on compilers supporting it? (Done 2012/09/15).
- The front-end should make the differences between implicit and explicit construction (Done 2012/09/15).
- The tutorial should add more examples concerning implicit or explicit conversions. (Done 2012/09/15).
- The documentation must explain how move semantics helps in this domain and what the backend needs to do to profit from this optimization. (Done 2012/09/15).
- The documentation should contain Throws specification on the `mp_number` and backend requirements operations. (Done 2012/09/15).
- The library interface should use the `noexcept` (`BOOST_NOEXCEPT, ...`) facilities (Done 2012/09/15).
- It is unfortunate that the generic `mp_number` front end can not make use `constexpr` as not all the backends can ensure this (done - we can go quite a way).
- literals: The library doesn't provide some kind of literals. I think that the `mp_number` class should provide a way to create literals if the backend is able to. (Done 2012/09/15).
- The `ExpressionTemplate` parameter could be defaulted to a traits class for more sensible defaults (done 2012/09/20).
- In `a = exp1 op exp2` where `a` occurs inside one of `exp1` or `exp2` then we can optimise and eliminate one more temporary (done 2012/09/20).

## Pre-Review Comments

- Make fixed precision orthogonal to Allocator type in `cpp_int`. Possible solution - add an additional `MaxBits` template argument that defaults to 0 (meaning keep going till no more space/memory). Done.
- Can ring types (exact floating point types) be supported? The answer should be yes, but someone needs to write it (Moved to TODO list).
- Should there be a choice of rounding mode (probably MPFR specific)? Moved to TODO list.
- Make the exponent type for `cpp_dec_float` a template parameter, maybe include support for big-integer exponents. Open question - what should be the default - `int32_t` or `int64_t`? (done 2012/09/06)
- Document the size requirements of fixed precision ints (done 2012/09/15).
- Document std lib function accuracy (done 2012/09/15).
- Be a bit clearer on the effects of sign-magnitude representation of `cpp_int` - `min == -max` etc - done.
- Document `cpp_dec_float` precision, rounding, and exponent size (done 2012/09/06).
- Can we be clearer in the docs that mixed arithmetic doesn't work (no longer applicable as of 2012/09/06)?
- Document round functions behaviour better (they behave as in C++11) (added note 2012/09/06).
- Document limits on size of `cpp_dec_float` (done 2012/09/06).
- Add support for fused multiply add (and subtract). GMP `mpz_t` could use this (done 2012/09/20).

## FAQ

Why do I get compiler errors when passing a `number` to a template function?

Most likely you are actually passing an expression template type to the function and template-argument-deduction deduces the "wrong" type. Try casting the arguments involving expressions

	to the actual number type, or as a last resort turning off expression template support in the number type you are using.
When is expression template support a performance gain?	As a general rule, expression template support adds a small runtime overhead creating and unpacking the expression templates, but greatly reduces the number of temporaries created. So it's most effective in improving performance when the cost of creating a temporary is high: for example when creating a temporary involves a memory allocation. It is least effective (and may even be a dis-optimisation) when temporaries are cheap: for example if the number type is basically a thin wrapper around a native arithmetic type. In addition, since the library makes extensive use of thin inline wrapper functions, turning on compiler optimization is essential to achieving high performance.
Do expression templates reorder operations?	Yes they do, sometimes quite radically so, if this is a concern then they should be turned off for the number type you are using.
I can't construct my number type from <i>some other type</i> , but the docs indicate that the conversion should be allowed, what's up?	Some conversions are <i>explicit</i> , that includes construction from a string, or constructing from any type that may result in loss of precision (for example constructing an integer type from a float).
Why do I get an exception thrown (or the program crash due to an uncaught exception) when using the bitwise operators on a checked <code>cpp_int</code> ?	Bitwise operations on negative values (or indeed any signed integer type) are unspecified by the standard. As a result any attempt to carry out a bitwise operation on a negative checked-integer will result in a <code>std::range_error</code> being thrown.
Why do I get compiler errors when trying to use the complement operator?	Use of the complement operator on signed types is problematic as the result is unspecified by the standard, and is further complicated by the fact that most extended precision integer types use a sign-magnitude representation rather than the 2's complement one favored by most native integer types. As a result the complement operator is deliberately disabled for checked <code>cpp_int</code> 's. Unchecked <code>cpp_int</code> 's give the same valued result as a 2's complement type would, but not the same bit-pattern.
Why can't I negate an unsigned type?	The unary negation operator is deliberately disabled for unsigned integer types as its use would almost always be a programming error.
Why doesn't the library use proto?	A very early version of the library did use proto, but compile times became too slow for the library to be usable. Since the library only required a tiny fraction of what proto has to offer anyway, a lightweight expression template mechanism was used instead. Compile times are still too slow...
Why not abstract out addition/multiplication algorithms?	This was deemed not to be practical: these algorithms are intimately tied to the actual data representation used.

## Acknowledgements

This library would not have happened without:

- Christopher Kormanyos' C++ decimal number code.
- Paul Bristow for patiently testing, and commenting on the library.
- All the folks at GMP, MPFR and libtommath, for providing the "guts" that makes this library work.
- "[The Art Of Computer Programming](#)", Donald E. Knuth, Volume 2: Seminumerical Algorithms, Third Edition (Reading, Massachusetts: Addison-Wesley, 1997), xiv+762pp. ISBN 0-201-89684-2

# Indexes

## Function Index

### A

- abs
  - number, 84
- add
  - Generic Integer Operations, 61
  - Mixed Precision Arithmetic, 59
  - number, 84
  - TODO, 150
- assign
  - number, 84
- assign\_components
  - Optional Requirements on the Backend Type, 103

### B

- bits
  - Generating Random Numbers, 53
  - Input Output, 81
  - Rounding Rules for Conversions, 58
  - std::numeric\_limits<> constants, 63
- bit\_flip
  - Generic Integer Operations, 61
  - number, 84
- bit\_set
  - Generic Integer Operations, 61
  - number, 84
- bit\_test
  - Generic Integer Operations, 61
  - number, 84
- bit\_unset
  - Generic Integer Operations, 61
  - number, 84

### C

- compare
  - number, 84

### D

- data
  - float128, 25
  - gmp\_float, 21
  - gmp\_int, 12
  - gmp\_rational, 41
  - mpfi\_float, 36
  - mpfr\_float, 23
- default\_precision
  - number, 84
- divide\_qr
  - Generic Integer Operations, 61
  - number, 84



**E**

empty

    mpfi\_float, 36

eval\_acos

    Optional Requirements on the Backend Type, 103

eval\_add

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_asin

    Optional Requirements on the Backend Type, 103

eval\_atan

    Optional Requirements on the Backend Type, 103

eval\_atan2

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_and

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_or

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_xor

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_bit\_flip

    Optional Requirements on the Backend Type, 103

eval\_bit\_set

    Optional Requirements on the Backend Type, 103

eval\_bit\_test

    Optional Requirements on the Backend Type, 103

eval\_bit\_unset

    Optional Requirements on the Backend Type, 103

eval\_ceil

    Compulsory Requirements on the Backend type., 103

eval\_complement

    Compulsory Requirements on the Backend type., 103

eval\_convert\_to

    Compulsory Requirements on the Backend type., 103

eval\_cos

    Optional Requirements on the Backend Type, 103

eval\_cosh

    Optional Requirements on the Backend Type, 103

eval\_decrement

    Optional Requirements on the Backend Type, 103

eval\_divide

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_eq

    Optional Requirements on the Backend Type, 103

eval\_exp

    Optional Requirements on the Backend Type, 103

eval\_fabs

    Optional Requirements on the Backend Type, 103

eval\_floor

    Compulsory Requirements on the Backend type., 103

eval\_fmod

    Optional Requirements on the Backend Type, 103

eval\_frexp

    Compulsory Requirements on the Backend type., 103

eval\_gcd  
Optional Requirements on the Backend Type, 103

eval\_get\_sign  
Optional Requirements on the Backend Type, 103

eval\_gt  
Optional Requirements on the Backend Type, 103

eval\_increment  
Optional Requirements on the Backend Type, 103

eval\_integer\_sqrt  
Optional Requirements on the Backend Type, 103

eval\_is\_zero  
Optional Requirements on the Backend Type, 103

eval\_lcm  
Optional Requirements on the Backend Type, 103

eval\_ldexp  
Compulsory Requirements on the Backend type., 103

eval\_left\_shift  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_log  
Optional Requirements on the Backend Type, 103

eval\_log10  
Optional Requirements on the Backend Type, 103

eval\_lt  
Optional Requirements on the Backend Type, 103

eval\_modulus  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_msb  
Optional Requirements on the Backend Type, 103

eval\_multiply  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_multiply\_add  
Optional Requirements on the Backend Type, 103

eval\_multiply\_subtract  
Optional Requirements on the Backend Type, 103

eval\_pow  
Optional Requirements on the Backend Type, 103

eval\_powm  
Optional Requirements on the Backend Type, 103

eval\_qr  
Optional Requirements on the Backend Type, 103

eval\_right\_shift  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_round  
Optional Requirements on the Backend Type, 103

eval\_sin  
Optional Requirements on the Backend Type, 103

eval\_sinh  
Optional Requirements on the Backend Type, 103

eval\_sqrt  
Compulsory Requirements on the Backend type., 103

eval\_subtract  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_tan

Optional Requirements on the Backend Type, 103  
eval\_tanh  
Optional Requirements on the Backend Type, 103  
eval\_trunc  
Optional Requirements on the Backend Type, 103

## F

fpclassify  
number, 84

## I

if  
Primality Testing, 55  
in  
mpfi\_float, 36  
infinity  
std::numeric\_limits<> functions, 70  
integer\_modulus  
Generic Integer Operations, 61  
number, 84  
iround  
number, 84  
isfinite  
number, 84  
isinf  
number, 84  
isnan  
number, 84  
isnormal  
number, 84  
itrunc  
number, 84

## L

llround  
number, 84  
lltrunc  
number, 84  
log\_postfix\_event  
logged\_adaptor, 44  
log\_prefix\_event  
logged\_adaptor, 44  
lround  
number, 84  
lsb  
Generic Integer Operations, 61  
number, 84  
ltrunc  
number, 84

## M

max  
std::numeric\_limits<> constants, 63  
std::numeric\_limits<> functions, 70  
miller\_rabin\_test  
Generic Integer Operations, 61  
number, 84

Primality Testing, 55

min

std::numeric\_limits<> functions, 70

msb

Generic Integer Operations, 61

number, 84

multiply

Generic Integer Operations, 61

Mixed Precision Arithmetic, 59

number, 84

## O

overlap

mpfi\_float, 36

## P

powm

Generic Integer Operations, 61

number, 84

precision

FAQ, 151

Generic Integer Operations, 61

Introduction, 3

number, 84

proper\_subset

mpfi\_float, 36

## R

r

Bit Operations, 16

round

number, 84

## S

sign

number, 84

singleton

mpfi\_float, 36

sqrt

Generic Integer Operations, 61

number, 84

std::numeric\_limits<> functions, 70

str

number, 84

subset

mpfi\_float, 36

subtract

Generic Integer Operations, 61

Mixed Precision Arithmetic, 59

number, 84

swap

number, 84

## T

trunc

number, 84

two  
std::numeric\_limits<> constants, 63

## V

value  
cpp\_bin\_float, 101  
number, 84  
std::numeric\_limits<> functions, 70

## Z

zero  
gmp\_float, 22  
gmp\_int, 13  
gmp\_rational, 42  
std::numeric\_limits<> constants, 69  
tommath\_rational, 43  
tom\_int, 14  
zero\_in  
mpfi\_float, 36

# Class Index

## C

component\_type  
number, 84  
cpp\_bin\_float  
cpp\_bin\_float, 18, 100  
cpp\_dec\_float  
cpp\_dec\_float, 20, 102  
cpp\_int\_backend  
cpp\_int, 9, 97

## D

debug\_adaptor  
debug\_adaptor, 46

## E

expression\_template\_default  
number, 84

## F

float128\_backend  
float128, 25

## G

gmp\_float  
gmp\_float, 21, 99  
gmp\_int  
gmp\_int, 12, 98  
gmp\_rational  
gmp\_rational, 41

## I

is\_explicitly\_convertible  
Internal Support Code, 102

is\_lossy\_conversion  
    Internal Support Code, 102  
is\_number  
    number, 84  
is\_number\_expression  
    number, 84  
is\_restricted\_conversion  
    Internal Support Code, 102  
is\_signed\_number  
    Internal Support Code, 102  
is\_unsigned\_number  
    Internal Support Code, 102

## L

logged\_adaptor  
    logged\_adaptor, 44

## M

mpfi\_float\_backend  
    mpfi\_float, 36  
mpfr\_float\_backend  
    mpfr\_float, 23, 100  
    mpfr\_float\_backend, 100

## N

number  
    number, 84  
number\_category  
    number, 84

## T

tommath\_int  
    tom\_int, 13, 99

# Typedef Index

## C

checked\_cpp\_int  
    cpp\_int, 9, 97  
checked\_cpp\_rational  
    cpp\_int, 9, 97  
checked\_cpp\_rational\_backend  
    cpp\_int, 9, 97  
checked\_int1024\_t  
    cpp\_int, 9, 97  
checked\_int128\_t  
    cpp\_int, 9, 97  
checked\_int256\_t  
    cpp\_int, 9, 97  
checked\_int512\_t  
    cpp\_int, 9, 97  
checked\_uint1024\_t  
    cpp\_int, 9, 97  
checked\_uint128\_t  
    cpp\_int, 9, 97  
checked\_uint256\_t

- cpp\_int, 9, 97
- checked\_uint512\_t
  - cpp\_int, 9, 97
- cpp\_bin\_float\_100
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_50
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_double
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_double\_extended
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_quad
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_single
  - cpp\_bin\_float, 18, 100
- cpp\_dec\_float\_100
  - cpp\_dec\_float, 20, 102
- cpp\_dec\_float\_50
  - cpp\_dec\_float, 20, 102
- std::numeric\_limits<> constants, 63
- std::numeric\_limits<> functions, 70
- cpp\_int
  - cpp\_int, 9, 97
- cpp\_rational
  - cpp\_int, 9, 97
  - cpp\_rational, 40
- cpp\_rational\_backend
  - cpp\_int, 9, 97
  - cpp\_rational, 40

## F

- float128
  - float128, 25
  - std::numeric\_limits<> constants, 63

## I

- int1024\_t
  - cpp\_int, 9, 97
- int128\_t
  - cpp\_int, 9, 97
- int256\_t
  - cpp\_int, 9, 97
- int512\_t
  - cpp\_int, 9, 97
- int\_type
  - Primality Testing, 55

## L

- limb\_type
  - cpp\_int, 9, 97

## M

- mpfi\_float
  - mpfi\_float, 36
- mpfi\_float\_1000
  - mpfi\_float, 36
- mpfi\_float\_50

- mpfi\_float, 36
- mpfr\_float
  - mpfr\_float, 23
  - mpfr\_float\_backend, 23, 100
- mpfr\_float\_100
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_1000
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_50
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_500
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpf\_float
  - gmp\_float, 21, 99
- mpf\_float\_100
  - gmp\_float, 21, 99
- mpf\_float\_1000
  - gmp\_float, 21, 99
- mpf\_float\_50
  - gmp\_float, 21, 99
- mpf\_float\_500
  - gmp\_float, 21, 99
- mpq\_rational
  - gmp\_rational, 41
- mpz\_int
  - gmp\_int, 12, 98
- mp\_type
  - Calculating an Integral, 32
  - Polynomial Evaluation, 34

## S

- static\_mpfr\_float\_100
  - mpfr\_float, 23
- static\_mpfr\_float\_50
  - mpfr\_float, 23

## T

- tommath\_rational
  - tommath\_rational, 42
- tom\_int
  - tom\_int, 13, 99
- tom\_rational
  - tommath\_rational, 42

## U

- uint1024\_t
  - cpp\_int, 9, 97
- uint128\_t
  - cpp\_int, 9, 97
- uint256\_t
  - cpp\_int, 9, 97
- uint512\_t
  - cpp\_int, 9, 97



# Index

## A

- abs
  - number, 84
- add
  - Generic Integer Operations, 61
  - Mixed Precision Arithmetic, 59
  - number, 84
  - TODO, 150
- assign
  - number, 84
- assign\_components
  - Optional Requirements on the Backend Type, 103

## B

- Bit Operations
  - r, 16
- bits
  - Generating Random Numbers, 53
  - Input Output, 81
  - Rounding Rules for Conversions, 58
  - std::numeric\_limits<> constants, 63
- bit\_flip
  - Generic Integer Operations, 61
  - number, 84
- bit\_set
  - Generic Integer Operations, 61
  - number, 84
- bit\_test
  - Generic Integer Operations, 61
  - number, 84
- bit\_unset
  - Generic Integer Operations, 61
  - number, 84
- BOOST\_MP\_DEFINE\_SIZED\_CPP\_INT\_LITERAL
  - Literal Types and constexpr Support, 56
- BOOST\_MP\_MIN\_EXPONENT\_DIGITS
  - Input Output, 81
- BOOST\_MP\_USE\_FLOAT128
  - float128, 26
- BOOST\_MP\_USE\_QUAD
  - float128, 26

## C

- Calculating an Integral
  - mp\_type, 32
- checked\_cpp\_int
  - cpp\_int, 9, 97
- checked\_cpp\_rational
  - cpp\_int, 9, 97
- checked\_cpp\_rational\_backend
  - cpp\_int, 9, 97
- checked\_int1024\_t
  - cpp\_int, 9, 97
- checked\_int128\_t

- cpp\_int, 9, 97
- checked\_int256\_t
  - cpp\_int, 9, 97
- checked\_int512\_t
  - cpp\_int, 9, 97
- checked\_uint1024\_t
  - cpp\_int, 9, 97
- checked\_uint128\_t
  - cpp\_int, 9, 97
- checked\_uint256\_t
  - cpp\_int, 9, 97
- checked\_uint512\_t
  - cpp\_int, 9, 97
- compare
  - number, 84
- component\_type
  - number, 84
- Compulsory Requirements on the Backend type.
  - eval\_add, 103
  - eval\_bitwise\_or, 103
  - eval\_bitwise\_xor, 103
  - eval\_ceil, 103
  - eval\_complement, 103
  - eval\_convert\_to, 103
  - eval\_divide, 103
  - eval\_floor, 103
  - eval\_frexp, 103
  - eval\_ldexp, 103
  - eval\_left\_shift, 103
  - eval\_modulus, 103
  - eval\_multiply, 103
  - eval\_right\_shift, 103
  - eval\_sqrt, 103
  - eval\_subtract, 103
- cpp\_bin\_float
  - cpp\_bin\_float, 18, 100
  - cpp\_bin\_float\_100, 18, 100
  - cpp\_bin\_float\_50, 18, 100
  - cpp\_bin\_float\_double, 18, 100
  - cpp\_bin\_float\_double\_extended, 18, 100
  - cpp\_bin\_float\_quad, 18, 100
  - cpp\_bin\_float\_single, 18, 100
  - value, 101
- cpp\_bin\_float\_100
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_50
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_double
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_double\_extended
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_quad
  - cpp\_bin\_float, 18, 100
- cpp\_bin\_float\_single
  - cpp\_bin\_float, 18, 100
- cpp\_dec\_float
  - cpp\_dec\_float, 20, 102
  - cpp\_dec\_float\_100, 20, 102

- cpp\_dec\_float\_50, 20, 102
- cpp\_dec\_float\_100
  - cpp\_dec\_float, 20, 102
- cpp\_dec\_float\_50
  - cpp\_dec\_float, 20, 102
  - std::numeric\_limits<> constants, 63
  - std::numeric\_limits<> functions, 70
- cpp\_int
  - checked\_cpp\_int, 9, 97
  - checked\_cpp\_rational, 9, 97
  - checked\_cpp\_rational\_backend, 9, 97
  - checked\_int1024\_t, 9, 97
  - checked\_int128\_t, 9, 97
  - checked\_int256\_t, 9, 97
  - checked\_int512\_t, 9, 97
  - checked\_uint1024\_t, 9, 97
  - checked\_uint128\_t, 9, 97
  - checked\_uint256\_t, 9, 97
  - checked\_uint512\_t, 9, 97
  - cpp\_int, 9, 97
  - cpp\_int\_backend, 9, 97
  - cpp\_rational, 9, 97
  - cpp\_rational\_backend, 9, 97
  - int1024\_t, 9, 97
  - int128\_t, 9, 97
  - int256\_t, 9, 97
  - int512\_t, 9, 97
  - limb\_type, 9, 97
  - uint1024\_t, 9, 97
  - uint128\_t, 9, 97
  - uint256\_t, 9, 97
  - uint512\_t, 9, 97
- cpp\_int\_backend
  - cpp\_int, 9, 97
- cpp\_rational
  - cpp\_int, 9, 97
  - cpp\_rational, 40
  - cpp\_rational\_backend, 40
- cpp\_rational\_backend
  - cpp\_int, 9, 97
  - cpp\_rational, 40

## D

- data
  - float128, 25
  - gmp\_float, 21
  - gmp\_int, 12
  - gmp\_rational, 41
  - mpfi\_float, 36
  - mpfr\_float, 23
- debug\_adaptor
  - debug\_adaptor, 46
- default\_precision
  - number, 84
- divide\_qr
  - Generic Integer Operations, 61
  - number, 84

**E**

empty

    mpfi\_float, 36

eval\_acos

    Optional Requirements on the Backend Type, 103

eval\_add

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_asin

    Optional Requirements on the Backend Type, 103

eval\_atan

    Optional Requirements on the Backend Type, 103

eval\_atan2

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_and

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_or

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_bitwise\_xor

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_bit\_flip

    Optional Requirements on the Backend Type, 103

eval\_bit\_set

    Optional Requirements on the Backend Type, 103

eval\_bit\_test

    Optional Requirements on the Backend Type, 103

eval\_bit\_unset

    Optional Requirements on the Backend Type, 103

eval\_ceil

    Compulsory Requirements on the Backend type., 103

eval\_complement

    Compulsory Requirements on the Backend type., 103

eval\_convert\_to

    Compulsory Requirements on the Backend type., 103

eval\_cos

    Optional Requirements on the Backend Type, 103

eval\_cosh

    Optional Requirements on the Backend Type, 103

eval\_decrement

    Optional Requirements on the Backend Type, 103

eval\_divide

    Compulsory Requirements on the Backend type., 103

    Optional Requirements on the Backend Type, 103

eval\_eq

    Optional Requirements on the Backend Type, 103

eval\_exp

    Optional Requirements on the Backend Type, 103

eval\_fabs

    Optional Requirements on the Backend Type, 103

eval\_floor

    Compulsory Requirements on the Backend type., 103

eval\_fmod

    Optional Requirements on the Backend Type, 103

eval\_frexp

    Compulsory Requirements on the Backend type., 103

eval\_gcd  
Optional Requirements on the Backend Type, 103

eval\_get\_sign  
Optional Requirements on the Backend Type, 103

eval\_gt  
Optional Requirements on the Backend Type, 103

eval\_increment  
Optional Requirements on the Backend Type, 103

eval\_integer\_sqrt  
Optional Requirements on the Backend Type, 103

eval\_is\_zero  
Optional Requirements on the Backend Type, 103

eval\_lcm  
Optional Requirements on the Backend Type, 103

eval\_ldexp  
Compulsory Requirements on the Backend type., 103

eval\_left\_shift  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_log  
Optional Requirements on the Backend Type, 103

eval\_log10  
Optional Requirements on the Backend Type, 103

eval\_lt  
Optional Requirements on the Backend Type, 103

eval\_modulus  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_msb  
Optional Requirements on the Backend Type, 103

eval\_multiply  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_multiply\_add  
Optional Requirements on the Backend Type, 103

eval\_multiply\_subtract  
Optional Requirements on the Backend Type, 103

eval\_pow  
Optional Requirements on the Backend Type, 103

eval\_powm  
Optional Requirements on the Backend Type, 103

eval\_qr  
Optional Requirements on the Backend Type, 103

eval\_right\_shift  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_round  
Optional Requirements on the Backend Type, 103

eval\_sin  
Optional Requirements on the Backend Type, 103

eval\_sinh  
Optional Requirements on the Backend Type, 103

eval\_sqrt  
Compulsory Requirements on the Backend type., 103

eval\_subtract  
Compulsory Requirements on the Backend type., 103  
Optional Requirements on the Backend Type, 103

eval\_tan

- Optional Requirements on the Backend Type, 103
- eval\_tanh
  - Optional Requirements on the Backend Type, 103
- eval\_trunc
  - Optional Requirements on the Backend Type, 103
- expression\_template\_default
  - number, 84

## F

### FAQ

- precision, 151
- float128
  - BOOST\_MP\_USE\_FLOAT128, 26
  - BOOST\_MP\_USE\_QUAD, 26
  - data, 25
  - float128, 25
  - float128\_backend, 25
  - std::numeric\_limits<> constants, 63
- float128\_backend
  - float128, 25
- fpclassify
  - number, 84

## G

### Generating Random Numbers

- bits, 53

### Generic Integer Operations

- add, 61
- bit\_flip, 61
- bit\_set, 61
- bit\_test, 61
- bit\_unset, 61
- divide\_qr, 61
- integer\_modulus, 61
- lsb, 61
- miller\_rabin\_test, 61
- msb, 61
- multiply, 61
- powm, 61
- precision, 61
- sqrt, 61
- subtract, 61
- gmp\_float
  - data, 21
  - gmp\_float, 21, 99
  - mpf\_float, 21, 99
  - mpf\_float\_100, 21, 99
  - mpf\_float\_1000, 21, 99
  - mpf\_float\_50, 21, 99
  - mpf\_float\_500, 21, 99
  - zero, 22
- gmp\_int
  - data, 12
  - gmp\_int, 12, 98
  - mpz\_int, 12, 98
  - zero, 13
- gmp\_rational

- data, 41
- gmp\_rational, 41
- mpq\_rational, 41
- zero, 42

## I

- if
  - Primality Testing, 55
- in
  - mpfi\_float, 36
- infinity
  - std::numeric\_limits<> functions, 70
- Input Output
  - bits, 81
  - BOOST\_MP\_MIN\_EXPONENT\_DIGITS, 81
- int1024\_t
  - cpp\_int, 9, 97
- int128\_t
  - cpp\_int, 9, 97
- int256\_t
  - cpp\_int, 9, 97
- int512\_t
  - cpp\_int, 9, 97
- integer\_modulus
  - Generic Integer Operations, 61
  - number, 84
- Internal Support Code
  - is\_explicitly\_convertible, 102
  - is\_lossy\_conversion, 102
  - is\_restricted\_conversion, 102
  - is\_signed\_number, 102
  - is\_unsigned\_number, 102
- Introduction
  - precision, 3
- int\_type
  - Primality Testing, 55
- iround
  - number, 84
- isfinite
  - number, 84
- isinf
  - number, 84
- isnan
  - number, 84
- isnormal
  - number, 84
- is\_explicitly\_convertible
  - Internal Support Code, 102
- is\_lossy\_conversion
  - Internal Support Code, 102
- is\_number
  - number, 84
- is\_number\_expression
  - number, 84
- is\_restricted\_conversion
  - Internal Support Code, 102
- is\_signed\_number

- Internal Support Code, 102
- is\_unsigned\_number
  - Internal Support Code, 102
- itrunc
  - number, 84

## L

- limb\_type
  - cpp\_int, 9, 97
- Literal Types and constexpr Support
  - BOOST\_MP\_DEFINE\_SIZED\_CPP\_INT\_LITERAL, 56
- llround
  - number, 84
- lltrunc
  - number, 84
- logged\_adaptor
  - logged\_adaptor, 44
  - log\_postfix\_event, 44
  - log\_prefix\_event, 44
- log\_postfix\_event
  - logged\_adaptor, 44
- log\_prefix\_event
  - logged\_adaptor, 44
- lround
  - number, 84
- lsb
  - Generic Integer Operations, 61
  - number, 84
- ltrunc
  - number, 84

## M

- max
  - std::numeric\_limits<> constants, 63
  - std::numeric\_limits<> functions, 70
- miller\_rabin\_test
  - Generic Integer Operations, 61
  - number, 84
  - Primality Testing, 55
- min
  - std::numeric\_limits<> functions, 70
- Mixed Precision Arithmetic
  - add, 59
  - multiply, 59
  - subtract, 59
- mpfi\_float
  - data, 36
  - empty, 36
  - in, 36
  - mpfi\_float, 36
  - mpfi\_float\_1000, 36
  - mpfi\_float\_50, 36
  - mpfi\_float\_backend, 36
  - overlap, 36
  - proper\_subset, 36
  - singleton, 36
  - subset, 36



- zero\_in, 36
- mpfi\_float\_1000
  - mpfi\_float, 36
- mpfi\_float\_50
  - mpfi\_float, 36
- mpfi\_float\_backend
  - mpfi\_float, 36
- mpfr\_float
  - data, 23
  - mpfr\_float, 23
  - mpfr\_float\_100, 23
  - mpfr\_float\_1000, 23
  - mpfr\_float\_50, 23
  - mpfr\_float\_500, 23
  - mpfr\_float\_backend, 23, 100
  - static\_mpfr\_float\_100, 23
  - static\_mpfr\_float\_50, 23
- mpfr\_float\_100
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_1000
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_50
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_500
  - mpfr\_float, 23
  - mpfr\_float\_backend, 100
- mpfr\_float\_backend
  - mpfr\_float, 23, 100
  - mpfr\_float\_100, 100
  - mpfr\_float\_1000, 100
  - mpfr\_float\_50, 100
  - mpfr\_float\_500, 100
  - mpfr\_float\_backend, 100
- mpf\_float
  - gmp\_float, 21, 99
- mpf\_float\_100
  - gmp\_float, 21, 99
- mpf\_float\_1000
  - gmp\_float, 21, 99
- mpf\_float\_50
  - gmp\_float, 21, 99
- mpf\_float\_500
  - gmp\_float, 21, 99
- mpq\_rational
  - gmp\_rational, 41
- mpz\_int
  - gmp\_int, 12, 98
- mp\_type
  - Calculating an Integral, 32
  - Polynomial Evaluation, 34
- msb
  - Generic Integer Operations, 61
- number, 84
- multiply
  - Generic Integer Operations, 61

Mixed Precision Arithmetic, 59  
number, 84

## N

### number

- abs, 84
- add, 84
- assign, 84
- bit\_flip, 84
- bit\_set, 84
- bit\_test, 84
- bit\_unset, 84
- compare, 84
- component\_type, 84
- default\_precision, 84
- divide\_qr, 84
- expression\_template\_default, 84
- fpclassify, 84
- integer\_modulus, 84
- iround, 84
- isfinite, 84
- isinf, 84
- isnan, 84
- isnormal, 84
- is\_number, 84
- is\_number\_expression, 84
- itrunc, 84
- llround, 84
- lltrunc, 84
- lround, 84
- lsb, 84
- ltrunc, 84
- millerrabin\_test, 84
- msb, 84
- multiply, 84
- number, 84
- number\_category, 84
- powm, 84
- precision, 84
- round, 84
- sign, 84
- sqr, 84
- str, 84
- subtract, 84
- swap, 84
- trunc, 84
- value, 84

### number\_category

- number, 84

## O

### Optional Requirements on the Backend Type

- assign\_components, 103
- eval\_acos, 103
- eval\_add, 103
- eval\_asin, 103
- eval\_atan, 103

- eval\_atan2, 103
- eval\_bitwise\_and, 103
- eval\_bitwise\_or, 103
- eval\_bitwise\_xor, 103
- eval\_bit\_flip, 103
- eval\_bit\_set, 103
- eval\_bit\_test, 103
- eval\_bit\_unset, 103
- eval\_cos, 103
- eval\_cosh, 103
- eval\_decrement, 103
- eval\_divide, 103
- eval\_eq, 103
- eval\_exp, 103
- eval\_fabs, 103
- eval\_fmod, 103
- eval\_gcd, 103
- eval\_get\_sign, 103
- eval\_gt, 103
- eval\_increment, 103
- eval\_integer\_sqrt, 103
- eval\_is\_zero, 103
- eval\_lcm, 103
- eval\_left\_shift, 103
- eval\_log, 103
- eval\_log10, 103
- eval\_lt, 103
- eval\_modulus, 103
- eval\_msb, 103
- eval\_multiply, 103
- eval\_multiply\_add, 103
- eval\_multiply\_subtract, 103
- eval\_pow, 103
- eval\_powm, 103
- eval\_qr, 103
- eval\_right\_shift, 103
- eval\_round, 103
- eval\_sin, 103
- eval\_sinh, 103
- eval\_subtract, 103
- eval\_tan, 103
- eval\_tanh, 103
- eval\_trunc, 103

overlap

- mpfi\_float, 36

## P

Polynomial Evaluation

- mp\_type, 34

powm

- Generic Integer Operations, 61
- number, 84

precision

- FAQ, 151
- Generic Integer Operations, 61
- Introduction, 3
- number, 84

## Primality Testing

- if, 55
  - int\_type, 55
  - miller\_rabin\_test, 55
- proper\_subset
- mpfi\_float, 36

**R**

- r
- Bit Operations, 16
- round
- number, 84
- Rounding Rules for Conversions
- bits, 58

**S**

- sign
- number, 84
- singleton
- mpfi\_float, 36
- sqrt
- Generic Integer Operations, 61
  - number, 84
  - std::numeric\_limits<> functions, 70
- static\_mpfr\_float\_100
- mpfr\_float, 23
- static\_mpfr\_float\_50
- mpfr\_float, 23
- std::numeric\_limits<> constants
- bits, 63
  - cpp\_dec\_float\_50, 63
  - float128, 63
  - max, 63
  - two, 63
  - zero, 69
- std::numeric\_limits<> functions
- cpp\_dec\_float\_50, 70
  - infinity, 70
  - max, 70
  - min, 70
  - sqrt, 70
  - value, 70
- str
- number, 84
- subset
- mpfi\_float, 36
- subtract
- Generic Integer Operations, 61
  - Mixed Precision Arithmetic, 59
  - number, 84
- swap
- number, 84

**T**

- TODO
- add, 150
- tommath\_int

- tom\_int, 13, 99
- tommath\_rational
  - tommath\_rational, 42
  - tom\_rational, 42
  - zero, 43
- tom\_int
  - tommath\_int, 13, 99
  - tom\_int, 13, 99
  - zero, 14
- tom\_rational
  - tommath\_rational, 42
- trunc
  - number, 84
- two
  - std::numeric\_limits<> constants, 63

## U

- uint1024\_t
  - cpp\_int, 9, 97
- uint128\_t
  - cpp\_int, 9, 97
- uint256\_t
  - cpp\_int, 9, 97
- uint512\_t
  - cpp\_int, 9, 97

## V

- value
  - cpp\_bin\_float, 101
  - number, 84
  - std::numeric\_limits<> functions, 70

## Z

- zero
  - gmp\_float, 22
  - gmp\_int, 13
  - gmp\_rational, 42
  - std::numeric\_limits<> constants, 69
  - tommath\_rational, 43
  - tom\_int, 14
- zero\_in
  - mpfi\_float, 36