

Document No: N4286

Supersedes: N4134

Date: 2014-11-18

Reply to: Gor Nishanov (gorn@microsoft.com), Jim Radigan (jradigan@microsoft.com)

Resumable Functions (revision 3)

Contents

Revisions and History	2
Terms and Definitions	3
Coroutine	3
Coroutine State / Coroutine Frame	3
Coroutine Promise	3
Coroutine Object / Coroutine Handle / Return Object of the Coroutine	3
Generator	3
Stackless Coroutine	3
Stackful Coroutine / Fiber / User-Mode thread	3
Split Stack / Linked Stack / Segmented Stack	3
Resumable Function	3
Discussion	4
Stackless vs Stackful	4
Implementation Experience	4
Asynchronous I/O	5
Generator	5
Parent-stealing parallel computations	6
Go-like channels and goroutines	7
Reactive Streams	7
Resumable lambdas as generator expressions	8
Conceptual Model	8
Resumable Function	8
Coroutine traits	9
Allocation and parameter copy optimizations	10
auto / decltype(auto) return type	10
C++ does not need generator expressions... it already has them!	10

Coroutine promise Requirements	11
Resumption function object	13
await operator	14
Evaluation of await expression	14
yield statement	16
Return statement.....	16
await-for statement	17
Trivial awaitable types	18
An expository Resumable Function Implementation	18
Coroutines in environments where exceptions are unavailable / banned.....	19
Allocation failure.....	20
Generalizing coroutine's promise set_exception	20
Await expression: Unwrapping of an eventual value	20
Await expression: Failure to launch an asynchronous operation.....	21
Asynchronous cancellation	21
Stateful Allocators Support.....	22
Override Selection of Coroutine Traits	23
Proposed Standard Wording.....	24
Acknowledgments.....	24
References	24
Appendix A: An example of generator coroutine implementation	25
Appendix B: boost::future adapters	27
Appendix C: Awaitable adapter over OS async facilities.....	28
Appendix D: Exceptionless error propagation with boost::future.....	29

Revisions and History

This document supersedes N4134. Changes relative to N4134 includes renaming resumable_traits and resumable_handle to coroutine_traits and coroutine_handle, adding a section describing a way to override default coroutine_traits selection, adjustment to coroutine_traits to support stateful allocators and altering Coroutine Promise concept requirements to allow omitting set_exception from the promise if stopping exception from propagating into a caller is not needed for a particular coroutine type.

Terms and Definitions

Coroutine

A generalized routine that in addition to traditional subroutine operations such as invoke and return supports suspend and resume operations.

Coroutine State / Coroutine Frame

A state that is created when coroutine is first invoked and destroyed once coroutine execution completes. Coroutine state includes a coroutine promise, formal parameters, variables and temporaries with automatic storage duration declared in the coroutine body and an implementation defined platform context. A platform context may contain room to save and restore platform specific data as needed to implement suspend and resume operations.

Coroutine Promise

A coroutine promise contains library specific data required for implementation of a higher-level abstraction exposed by a coroutine. For example, a coroutine implementing a task-like semantics providing an eventual value via `std::future<T>` is likely to have a coroutine promise that contains `std::promise<T>`. A coroutine implementing a generator may have a promise that stores a current value to be yielded and a state of the generator (active/cancelling/closed).

Coroutine Object / Coroutine Handle / Return Object of the Coroutine

An object returned from an initial invocation of a coroutine. A library developer defines the higher-level semantics exposed by the coroutine object. For example, generator coroutines can provide an input iterator that allows to consume values produced by the generator. For task-like coroutines, coroutine object can be used to obtain an eventual value (`future<T>`, for example).

Generator

A coroutine that provides a sequence of values. The body of the generator coroutine uses a **yield** statement to specify a value to be passed to the consumer. Emitting a value suspends the coroutine, invoking a pull operation on a channel resumes the coroutine.

Stackless Coroutine

A stackless coroutine is a coroutine which state includes variables and temporaries with automatic storage duration in the body of the coroutine and **does not** include the call stack.

Stackful Coroutine / Fiber / User-Mode thread

A stackful coroutine state **includes** the full call stack associated with its execution enabling suspension from nested stack frames. Stackful coroutines are equivalent to fibers or user-mode threads.

Split Stack / Linked Stack / Segmented Stack

A compiler / linker technology that enables non-contiguous stacks.

Resumable Function

Proposed C++ language mechanism to implement stackless coroutines.

Discussion

Motivation for extending C++ language and libraries to support coroutines was covered by papers N3858 (resumable functions) and N3985 (A proposal to add coroutines to C++ standard library) and will not be repeated here.

Design goals for this revision of resumable functions were to extend C++ language and standard library to support coroutines with the following characteristics:

- Highly scalable (to billions of concurrent coroutines).
- Highly efficient resume and suspend operations comparable in cost to a function call overhead.
- Seamless interaction with existing facilities with no overhead.
- Open ended coroutine machinery allowing library designers to develop coroutine libraries exposing various high-level semantics, such as generators, goroutines, tasks and more.
- Usable in environments where exception are forbidden or not available

Unlike N3985 (A proposal to add coroutine to the C++ standard library), which proposes certain high-level abstractions (coroutine-based input / output iterators), this paper focuses only on providing efficient language supported mechanism to suspend and resume a coroutine and leaves high-level semantics of what coroutines are to the discretion of a library developer and thus is comparable to Boost.Context rather than Boost.Coroutine / N3985.

Stackless vs Stackful

Design goals of scalability and seamless interaction with existing facilities without overhead (namely calling into existing libraries and OS APIs without restrictions) necessitates stackless coroutines.

General purpose stackful coroutines that reserve default stack for every coroutine (1MB on Windows, 2MB on Linux) will exhaust all available virtual memory in 32-bit address space with only a few thousand coroutines. Besides consuming virtual memory, stackful coroutines lead to memory fragmentation, since with common stack implementations, besides reserving virtual memory, the platform also commits first two pages of the stack (one as a read/write access to be used as a stack, another to act as a guard page to implement automatic stack growth), even though the actual state required by a coroutine could be as small as a few bytes.

A mitigation approach such as using split-stacks requires the entire program (including all the libraries and OS facilities it calls) to be either compiled with split-stacks or to incur run-time penalties when invoking code that is not compiled with split-stack support.

A mitigation approach such as using a small fixed sized stack limits what can be called from such coroutines as it must be guaranteed that none of the functions called shall ever consume more memory than allotted in a small fixed sized stack.

Implementation Experience

We implemented language changes described in this paper in Microsoft C++ compiler to gain experience and validate coroutine customization machinery. The following are illustrations of what library designers can achieve using coroutine mechanism described in this paper.

Note the usage of proposed, **await** operator, **yield** and *await-for* statements.

Asynchronous I/O

The following code implements zero-overhead abstractions over asynchronous socket API and windows threadpool.

```
std::future<void> tcp_reader(int total)
{
    char buf[64 * 1024];
    auto conn = await Tcp::Connect("127.0.0.1", 1337);
    do
    {
        auto bytesRead = await conn.read(buf, sizeof(buf));
        total -= bytesRead;
    }
    while (total > 0);
}

int main() { tcp_reader(1000 * 1000 * 1000).get(); }
```

Execution of this program incurs only one memory allocation¹² and no virtual function calls. The generated code is as good as or better than what could be written in C over raw OS facilities. The better part is due to the fact that OVERLAPPED structures (used in the implementation of `Tcp::Connect` and `conn.read`) are temporary objects on the frame of the coroutine whereas in traditional asynchronous C programs OVERLAPPED structures are dynamically allocated for every I/O operation (or for every distinct kind of I/O operation) on the heap.

Allocation of a future shared state (`N3936/[futures.state]`) associated with the future is combined with coroutine frame allocation and does not incur an extra allocation.

Generator

Another coroutine type was implemented to validate the generator pattern and a coroutine cancellation mechanics:

```
generator<int> fib(int n)
{
    int a = 0;
    int b = 1;
    while (n-- > 0)
    {
        yield a;
        auto next = a + b;
        a = b;
        b = next;
    }
}

int main()
{
    for (auto v : fib(35)) {
        std::cout << v << std::endl;
        if (v > 10)
            break;
    }
}
```

¹ Allocation of the frame of a resumable function

² not counting memory allocations incurred by OS facilities to perform an I/O

```

        break;
    }
}

```

Recursive application of generators allows to mitigate stackless coroutine inability to suspend from nested stack frames. This example is probably the most convoluted way to print number in range [1..100).

```

recursive_generator<int> range(int a, int b)
{
    auto n = b - a;

    if (n <= 0)
        return;

    if (n == 1)
    {
        yield a;
        return;
    }

    auto mid = a + n / 2;

    yield range(a, mid);
    yield range(mid, b);
}

int main() {
    for (auto v : range(1, 100))
        std::cout << v << std::endl;
}

```

Parent-stealing parallel computations

It is possible to adopt coroutine mechanics to support parallel scheduling techniques such as parent stealing [N3872].

```

spawnable<int> fib(int n) {
    if (n < 2) return n;
    return await(fib(n - 1) + fib(n - 2));
}

int main() { std::cout << fib(5).get() << std::endl; }

```

In this example **operator+** is overloaded for `spawnable<T>` type. Operator `+` schedules `fib(n - 2)` to a work queue whereas the execution continues into `fib(n-1)`. When eventual values for both `fib(n-1)` and `fib(n-2)` are ready, `fib(n)` is resumed and result of `await` expression is computed as the sum of the eventual value of the left and right operand to `+` operator.

Utilizing parent-stealing scheduling allows to compute `fib(42)` in less than 12k of space, whereas attempting to use more traditional scheduling will cause state explosion that will consume more than 2gig of memory around `fib(32)`.

Note, there are much better³ ways to compute Fibonacci numbers.

Go-like channels and goroutines

The following example (inspired by programming language go [GoLang]) creates one million goroutines, connects them to each other using channels and passes a value that will travel through all the coroutines.

```
goroutine pusher(channel<int>& left, channel<int>& right)
{
    for (;;) {
        auto val = await left.pull();
        await right.push(val + 1);
    }
}

int main()
{
    static const int N = 1000 * 1000;
    std::vector<channel<int>> c(N + 1);

    for (int i = 0; i < N; ++i)
        goroutine::go(pusher(c[i], c[i + 1]));

    c.front().sync_push(0);

    std::cout << c.back().sync_pull() << std::endl;
}
```

Reactive Streams

Resumable functions can be used as producers, consumers and transformers of reactive streams in recently rediscovered [Rx, ReactiveX, RxAtNetflix] functional reactive programming [FRP].

As a consumer (utilizing await-for statement proposed by this paper):

```
future<int> Sum(async_read_stream<int> & input)
{
    int sum = 0;
    for await (v : input)
    {
        sum += v;
    }
    return sum;
}
```

As a producer:

```
async_generator<int> Ticks()
{
    for(int tick = 0;; ++tick)
    {
        yield tick;
        await sleep_for(1ms);
    }
}
```

³ In constant space with only log n iterations

```
    }
}
```

As a transformer: (adds a timestamp to every observed value)

```
template<class T>
async_generator<pair<T, system_clock::time_point>>
Timestamp(async_read_stream<T> S)
{
    for await(v: S) yield {v, system_clock::now()};
}
```

Resumable lambdas as generator expressions

Resumable lambdas can be used as generator expressions [PythonGeneratorExpressions].

```
squares = (x*x for x in S) // python
auto squares = [&]{ for(x:S) yield x*x; } ; // C++
```

In this case squares is a lazy transformer of sequence S and similar in that respect to boost range adapters [BoostRangeAdapter].

Conceptual Model

Resumable Function

A function or a lambda is called **resumable** function or resumable lambda if a body of the function or lambda contains at least one suspend/resume point. Suspend/resume points are expressions with one or more *await operators*, *yield statements* or *await-for statements*. From this point on, we will use the term resumable function to refer to either resumable lambda or resumable function.

Suspend/resume points indicate the location where execution of the resumable function can be suspended and control returned to the current caller with an ability to resume execution at suspend/resume point later.

N3936/[intro.execution]/7 defines that suspension of a block preserves the automatic variables in a case of a function call or receipt of a signal. We propose to extend that language to coroutine suspension as well.

From the perspective of the caller, resumable function is just a normal function with that particular signature. The fact that a function is implemented as resumable function is unobservable by the caller. In fact, v1 version of some library can ship an implementation of some functions as resumable and switch it later to regular functions or vice versa without breaking any library user.

Design Note: Original design relied on resumable keyword to annotate resumable functions or lambdas. This proposal does away with resumable keyword relying on the presence of suspend/resume points. There were several motivations for this change.

1. It eliminates questions such as: Is resumable a part of signature or not? Does it change a calling conventions? Should it be specified only on a function definition?

2. It eliminates compilation errors due to the absence of resumable keyword that were in the category: “A compiler knows exactly what you mean, but won’t accept the code until you type ‘resumable’.”
3. Usability experience with the resumable functions implemented in C++ compiler by the authors. Initial implementation had resumable keyword and writing code felt unnecessarily verbose with having to type resumable in the declarations and definitions of functions or lambda expressions.

Coroutine traits

Coroutine traits are specialized by resumable functions to select an allocator and a coroutine promise to use in a particular resumable function.

If the signature of a resumable function is

R func(T1, T2, ... Tn)

then, a traits specialization `std::coroutine_traits<R,T1,T2,...,Tn>` will indicate what allocator and what coroutine promise to use.

For example, for coroutines returning `future<R>`, the following trait specialization can be provided.

```
template <typename R, typename... Ts>
struct coroutine_traits<std::future<R>, Ts...>
{
    template <typename... Us> static auto get_allocator(Us&&...);
    using promise_type = some-type-satisfying-coroutine-promise-requirements;
};
```

`get_allocator` should return an object satisfying allocator requirements

N3936/[allocator.requirements]. If `get_allocator` is not specified, a resumable function will default to using `std::allocator<char>`. `get_allocator` can examine the parameters and decide if there is a stateful allocator passed to a function and use it, otherwise, it can ignore the parameters and return a stateless allocator.

`promise_type` should satisfy requirements specified in “Coroutine Promise Requirements”. If

`promise_type` is not specified it is assumed to be as if defined as follows:

```
using promise_type = typename R::promise_type;
```

C++ standard library defines the coroutine traits as follows:

```
template <typename R, typename... Ts>
struct coroutine_traits
{
    template <typename... Us>
    static auto get_allocator(Us&&...) { return std::allocator<char>{}; }
    using promise_type = typename R::promise_type;
};
```

Design note: Another design option is to use only the return type in specializing the coroutine traits. The intention for including parameter types is to enable using parameters to alter allocation strategies or other implementation details while retaining the same coroutine return type.

Allocation and parameter copy optimizations

An invocation of a coroutine may incur an extra copy or move operation for the formal parameters if they need to be transferred from an ABI prescribed location into a memory allocated for the coroutine frame.

A parameter copy is not required if a coroutine never suspends or if it suspends but its parameters will not be accessed after the coroutine is resumed.

If a parameter copy/move is required, class object moves are performed according to the rules described in Copying and moving class objects section of the working draft standard 3936/[class.copy].

An implementation is allowed to elide calls to the allocator's allocate and deallocate functions and use stack memory of the caller instead if the meaning of the program will be unchanged except for the execution of the allocate and deallocate functions.

auto / decltype(auto) return type

If a function return type is auto or decltype(auto) and no trailing return type is specified, then the return type of the resumable function is deduced as follows:

1. If a yield statement and either an await expression or an await-for statement are present, then the return type is *default-async-generator*<T,R>, where T is deduced from the yield statements and R is deduced from return statements according to the rules of return type deduction described in N3936/[dcl.spec.auto].
2. Otherwise, if an await expression or an await-for statement are present in a function, then return type is *default-standard-task-type*<T> where type T is deduced from return statements as described in N3936/[dcl.spec.auto].
3. Otherwise, If a yield statement is present in a function, then return type is *default-generator-type*<T>, where T is deduced from the yield statements according to the rules of return type deduction described in N3936/[dcl.spec.auto].

At the moment we do not have a proposal for what *default-standard-task-type*, *default-generator-type* or *default-async-generator* should be. We envision that once resumable functions are available as a language feature, C++ community will come up with ingenious libraries utilizing that feature and some of them will get standardized and become *default-generator-type*, *default-task-type* and *default-async-generator* types. Appendix A, provides a sample implementation of generator<T>.

Until that time, an attempt to define resumable functions with auto / decltype(auto) and no trailing return type should result in a diagnostic message.

C++ does not need generator expressions... it already has them!

Assuming that we have a standard generator type that can be deduced as described before, the Python's generator expression can be trivially written in C++:

```
squares = (x*x for x in s) // python
auto squares = [&]{ for (x:s) yield x*x;} ; // C++
```

Coroutine promise Requirements

A library developer supplies the definition of the coroutine promise to implement desired high-level semantics associated with a coroutine type. The following tables describe the requirements on coroutine promise types.

Table 1: Descriptive Variable definitions

Variable	Definition
P	Coroutine promise type
P	A value of type P
E	A value of <code>std::exception_ptr</code> type
Rh	A value of type <code>std::coroutine_handle<P></code>
T	An arbitrary type T
V	A value of type T

Table 2: Coroutine Promise Requirements

Expression	Note
P{}	Constructs a promise type.
p.get_return_object(rh)	<p><code>get_return_object</code> is invoked by the coroutine to construct the return object prior to the first suspend operation.</p> <p><code>get_return_object</code> receives a value <code>std::coroutine_handle<Promise></code> as the first parameter.</p> <p>An object of <code>std::coroutine_handle<Promise></code> type can be used to resume the coroutine or get access to its promise.</p>
p.set_result(v)	<p>Sets the value associated with the promise. <code>set_result</code> is invoked by a resumable function when <code>return <expr>;</code> statement is encountered in a resumable function.</p> <p>If <code>p.set_result(v)</code> member function is not present, coroutine does not support eventual return value and presence of <code>return <expr></code> statement in the body is a compile time error.</p>
p.set_result()	<p><code>set_result()</code> is invoked by the resumable function when <code>return;</code> statement is encountered or the control reaches the end of the resumable function.</p> <p>If <code>set_result()</code> is not present, it is assumed that the function supports eventual value and diagnostic should be given if <code>return <expr></code> is not present in the body of the resumable function</p>

	<p>If both <code>set_result()</code> and <code>set_result(v)</code> are present in a promise type, the type does not satisfy coroutine promise requirement and a diagnostic message should be given to the user.</p>
<code>p.set_exception(e)</code>	<p><code>set_exception</code> is invoked by a resumable function when an unhandled exception occurs within a body of the resumable function.</p> <p>If promise does not provide <code>set_exception</code>, unhandled exceptions will propagate from a resumable functions normally.</p>
<code>p.yield_value(v)</code>	<p>returns: awaitable expression</p> <p><code>yield_value</code> is invoked when <code>yield</code> statement is evaluated in the resumable function.</p> <p>If <code>yield_value</code> member function is not present, using <code>yield</code> statement in the body of the resumable function results in a compile time error.</p> <p><code>yield <expr>;</code> is equivalent to <code>(void)(await <Promise>.yield_value(<expr>));</code></p> <p>Where a <code><Promise></code> refers to the coroutine promise of the enclosing resumable function.</p>
<code>initial_suspend()</code>	<p>Returns: awaitable expression</p> <p>A resumable function awaits on a value returned by the <code>initial_suspend()</code> member function immediately before user provided body of the resumable function is entered.</p> <p>This member function gives a library designer an option to suspend a coroutine after the coroutine frame is allocated, parameters are copied and return object is obtained, but before entering the user-provided body of the coroutine.</p> <p>For example, in a generator scenario, a library designer can choose to suspend a generator prior to invoking user provided body of the coroutine and to resume it once the user of the generator tries to pull the first value.</p>
<code>final_suspend()</code>	<p>Returns: nothrow awaitable expression Throws: nothing</p> <p>Resumable function awaits on a value returned by <code>final_suspend()</code> immediately after the user provided body of the resumable function i.e. point prior to the destruction and deallocation of a coroutine frame.</p>

	<p>This allows library designer to store the eventual value of the task, or the current value of the generator within the coroutine promise.</p> <p>Once the eventual value or last value is consumed, coroutine can be resumed to free up resources associated with it.</p>
cancellation_requested()	<p>Returns: bool to indicate whether coroutine is being cancelled</p> <p>cancellation_requested() is evaluated on resume code path. If it evaluates to true, control is transferred to the point immediately prior to compiler synthesized await promise-expr.final_suspend(), otherwise control is transferred to the current resume point.</p> <p>All of the objects with non-trivial destructors, will be destroyed in the same manner as if “goto end-label” statement was executed immediately after the resume point.</p> <p>(Assuming that “goto” was allowed to be used within an expression)</p>

Bikeshed: on_complete, on_error, on_next as a replacement for set_result, set_exception and yield_value, set_error as a replacement for set_exception.

Resumption function object

A resumable function has the ability to suspend evaluation by means of await operator or yield and await-for statements in its body. Evaluation may later be resumed at the suspend/resume point by invoking a resumption function object.

A resumption function object is defined by C++ standard library as follows:

```

template <typename Promise = void>
struct coroutine_handle;

template <> struct coroutine_handle<void>
{
    void operator() ();
    static coroutine_handle<void> from_address(void*);
    void * to_address() ;

    explicit operator bool() const;

    coroutine_handle() = default;
    explicit coroutine_handle(std::nullptr_t);
    coroutine_handle& operator = (nullptr_t);
};

template <typename Promise>
struct coroutine_handle: public coroutine_handle<>
{
    Promise & promise();
    Promise const & promise() const;
    static coroutine_handle<Promise> from_promise(Promise*);

```

```

        using coroutine_handle<>::coroutine_handle;
        coroutine_handle& operator = (nullptr_t);
};

template <typename Promise>
bool operator == (coroutine_handle<Promise> const& a,
                 coroutine_handle<Promise> const& b);

template <typename Promise>
bool operator != (coroutine_handle<Promise> const& a,
                 coroutine_handle<Promise> const& b)

```

Note, that by design, a resumption function object can be “round tripped” to void * and back. This property allows seamless interactions of resumable functions with existing C APIs⁴.

Resumption function object has two forms. One that provides an ability to resume evaluation of a resumable function and another, which additionally allows access to the coroutine promise of a particular resumable function.

Bikeshed: `resumption_handle`, `resumption_object`, `resumable_ptr`, `basic_coroutine_handle` instead of `coroutine_handle<void>`, `from_raw_address`, `to_raw_address`, `from_pvoid`, `to_pvoid`.

await operator

is a unary operator expression of the form: **await** *cast-expression*

1. The await operator shall not be invoked in a catch block of a try-statement⁵
2. The result of await is of type T, where T is the return type of the `await_resume` function invoked as described in the evaluation of await expression section. If T is void, then the await expression cannot be the operand of another expression.

Evaluation of await expression

An await expression in a form **await** *cast-expression* is equivalent to (if it were possible to write an expression in terms of a block, where return from the block becomes the result of the expression)

```

{
    auto && __expr = cast-expression;
    if ( !await-ready-expr ) {
        await-suspend-expr;
        suspend-resume-point
    }
    cancel-check;
    return await-resume-expr;
}

```

⁴ Usually C APIs take a callback and void* context. When the library/OS calls back into the user code, it invokes the callback passing the context back. `from_address()` static member function allows to reconstitute `coroutine_handle<>` from void* context and resume the coroutine

⁵ The motivation for this is to avoid interfering with existing exception propagation mechanisms, as they may be significantly (and negatively so) impacted should await be allowed to occur within exception handlers.

Where `__expr` is a variable defined for exposition only, and `_ExprT` is the type of the *cast-expression*, and `_PromiseT` is a type of the coroutine promise associated with current resumable function and the rest defined as follows:

await-ready-expr await-suspend-expr await-resume-expr	<p>— if <code>_ExprT</code> is a class type, the unqualified-ids await_ready, await_suspend and await_resume are looked up in the scope of class <code>_ExprT</code> as if by class member access lookup (N3936/3.4.5 [basic.lookup.classref]), and if it finds at least one declaration, then <code>await_ready-expr</code>, <code>await_suspend-expr</code> and <code>await_resume-expr</code> are <code>__expr.await_ready()</code>, <code>__expr.await_suspend(resumption-function-object)</code> and <code>__expr.await_resume()</code>, respectively;</p> <p>— otherwise, <code>await_ready-expr</code>, <code>await_suspend-expr</code> and <code>await_resume-expr</code> are <code>await_ready(__expr)</code>, <code>await_suspend(__expr, resumption-function-object)</code> and <code>await_resume(__expr)</code>, respectively, where <code>await_ready</code>, <code>await_suspend</code> and <code>await_resume</code> are looked up in the associated namespaces (N3936/3.4.2). [Note: Ordinary unqualified lookup (3.4.1) is not performed. —end note]</p> <p>A type for which <code>await_ready</code>, <code>await_suspend</code> and <code>await_resume</code> function can be looked up by the rules described above is called an awaitable type.</p> <p>If none of <code>await_xxx</code> functions can throw an exception, the awaitable type is called a nothrow awaitable type and expression of this type a nothrow awaitable expressions.</p>
resumption-function-object	A function object of type <code>std::coroutine_handle<_PromiseT></code> . When function object is invoked it will resume execution of the resumable function at the point marked by suspend-resume-point.
suspend-resume-point	When this point is reached, the coroutine is suspended. Once resumed, execution continues immediately after the suspend-resume-point
cancel-check	<p>For all await expressions except for the one implicitly synthesized by a compiler at the end of the function it is</p> <pre>if (<promise-expr>.cancellation_requested()) goto <end-label>;</pre> <p>where <code><promise-expr></code> is a reference to a coroutine promise associated with the current resumable function and an <code>end-label</code> is a label at the end of the user provided body of the resumable function, just prior to the <code>await <promise-expr>.final_suspend()</code>.</p>

Design Note: rules for lookup of `await_xxx` identifiers mirror the look up rules for range-based for statement. We also considered two other alternatives (we implemented all three approaches to test their usability, but found the other two less convenient than the one described above):

1. To have only ADL based lookup and not check for member functions. This approach was rejected as it disallowed one of the convenient patterns that was developed utilizing `await`. Namely to have compact declaration for asynchronous functions in a form: `auto Socket::AsyncRead(int count) { struct awaiter {...}; return awaiter{this, count}};`
2. Another approach considered and rejected was to have an **operator `await`** function found via ADL and having it to return an `awaitable_type` that should have `await_xxx` member functions defined. It was found more verbose than the proposed alternative.

Bikeshed: `await_suspend_needed`, `await_pre_suspend`, `await_pre_resume`

yield statement

A yield statement is a statement of form:

yield *expression*;

or

yield *braced-init-list*;

`yield` `<something>`; is equivalent to

`(void)(await <Promise>.yield_value(<something>))`

Where a `<Promise>` refers to the coroutine promise of the enclosing resumable function.

Design note: `yield` is a popular identifier, it is used in the standard library, e.g. `this_thread::yield()`. Introducing a `yield` keyword will break existing code. Having a two word keyword, such as **yield return** could be another choice.

Another alternative is to make **yield** a context-sensitive keyword that acts like a keyword at a statement level and as an identifier otherwise. To disambiguate access to a yield variable or a function, yield has to be prefixed by `::yield`, `->yield` and `.yield`. This will still break some existing code, but allows an escape hatch to patch up the code without having to rename yield function which could be defined by the libraries a user have no source access to.

Return statement

A return statement in a resumable function in a form **return** *expression*; is equivalent to:

`{ promise-expr.set_result(expression); goto end-label; }`

A return statement in a resumable function in a form **return** *braced-init-list*; is equivalent to:

`{ promise-expr.set_result(braced-init-list); goto end-label; }`

A return statement in a resumable function in a form **return**; is equivalent to:

`{ promise-expr.set_result(); goto end-label; }`

Where `end-label` is a label at the end of the user provided body of the resumable function, just prior to the `await <promise-expr>.final_suspend()`.

If resumable function does not have return statements in the form **return** *expression*; or **return** *braced-init-list*; then the function acts as if there is an implicit **return**; statement at the end of the function.

await-for statement

An await-for statement of the form:

```
for await ( for-range-declaration : expression ) statement
```

is equivalent to

```
{  
    auto && __range = expression;  
    for ( auto __begin = await begin-expr,  
        __end = end-expr;  
        __begin != __end;  
        await ++__begin )  
    {  
        for-range-declaration = *__begin;  
        statement  
    }  
}
```

where begin-expr and end-expr are defined as described in N3936/[stmt.ranged]/1.

The rationale for annotating begin-expr and ++ with await is as follows:

A model for consuming values from an asynchronous input stream looks like this:

```
for (;;) {  
    initiate async pull()  
    wait for completion of async  
    if (no more)  
        break;  
  
    process value  
}
```

We need to map this to iterators. The closest thing is an input_iterator.

For an input_iterator, frequent implementation of end() is a tag value that makes iterator equality comparison check for the end of the sequence, therefore, != end() is essentially an end-of-stream check.

So, begin() => async pull, therefore await is needed

++__begin => async pull, therefore await is needed

!= end() – is end-of-stream check post async pull, no need for await

Design note: We implemented two variants of *await-for* statement to evaluate their aesthetical appeal and typing convenience. One form was **for await**(x:S) , another **await for**(x:S)

Even though our initial choice was await for, we noticed that the brain was so hardwired to read things starting with **for** as loops, that **await for** did not register as loop when reading the code.

Trivial awaitable types

Standard library provides three awaitable types defined as follows:

```
namespace std {
    struct suspend_always {
        bool await_ready() const { return false; }
        void await_suspend(std::coroutine_handle<>) {}
        void await_resume() {}
    };
    struct suspend_never {
        bool await_ready() const { return true; }
        void await_suspend(std::coroutine_handle<>) {}
        void await_resume() {}
    };
    struct suspend_if
    {
        bool ready;

        suspend_if(bool condition): ready(!condition){}

        bool await_ready() const { return ready; }
        void await_suspend(std::coroutine_handle<>) {}
        void await_resume() {}
    };
}
```

These types are used in implementations of coroutine promises. Though they are trivial to implement, including them in the standard library eliminates the need for every library designer from doing their own implementation.

For example, `generator<T>` coroutine listed in the Appendix A, defines `yield_value` member function as follows:

```
std::suspend_always promise_type::yield_value(T const& v) {
    this->current_value = &v;
    return{};
}
```

An expository Resumable Function Implementation

Note: The following section is for illustration purposes only. It does not prescribe how resumable functions must be implemented.

Given a user authored function:

```
R foo(T1 a, T2 b) { body-containing-suspend-resume-points }
```

Compiler can constructs a function that behaves as if the following code was generated:

```
R foo(T1 a, T2 b) {
    using __traits = std::coroutine_traits<R, T1, T2>;
    struct __Context {
        __traits::promise_type _Promise;
        T1 a;
        T2 b;

        template <typename U1, typename U2>
```

```

__Context(U1&& a, U2&& b) : a(forward<U1>(a)), b(forward<U2>(b)) {}

void operator()() noexcept {
    await _Promise.initial_suspend();
    try { body-containing-suspend-resume-points-with-some-changes }
    catch (...) { _Promise.set_exception(std::current_exception()); }
__return_label:
    await _Promise.final_suspend();
    <deallocate-frame> (this, sizeof(__Context) + <X>);
}
};

auto mem = <allocate-frame>(sizeof(__Context) + <X>);
__Context * coro = nullptr;
try {
    coro = new (mem) __Context(a, b);
    auto result = __traits::get_return_object(
        std::coroutine_handle<__traits::promise_type>::from_promise(&coro->__Promise);
        (*coro)());
    return result;
}
catch (...) {
    if (coro) coro->~__Context();
    <deallocate-frame> (mem, sizeof(__Context) + <X>);
    throw;
}
}

```

Where, <X> is a constant representing the number of bytes that needs to be allocated to accommodate variables with automatic storage duration in the body of the resumable function and platform specific data that is needed to support resume and suspend.

Access to variables and temporaries with automatic storage duration in the body of operator() should be relative to “this” pointer at the offset equal to sizeof(*this).

<body-containing-suspend-resume-points-with-some-changes> is identical to <body-containing-suspend-resume-points> with the exception that await operators and yield, await-for statements, and return statements are transformed as described in earlier sections, __return_label is the <end-label> that return and await refer to and accesses to formal parameters a, b are replaced with accesses to member variables a and b of __Context class.

Coroutines in environments where exceptions are unavailable / banned

C++ exceptions represent a barrier to adoption of full power of C++. While this is unfortunate and may be rectified in the future, the current experience shows that kernel mode software, embedded software for devices and airplanes [JSF] forgo the use of C++ exceptions for various reasons.

Making coroutine fully dependent on C++ exceptions will limit their usefulness in contexts where asynchronous programming help is especially valuable (kernel mode drivers, embedded software, etc).

The following sections described how exceptions can be avoided in implementation and applications of resumable functions.

Allocation failure

To enable handling of allocation failures without relying on exception mechanism, `coroutine_traits` specialization can declare an optional static member function `get_return_object_on_allocation_failure`.

If `get_return_object_on_allocation_failure` member function is present, it is assumed that an allocator's `allocate` function will violate the standard requirements and will return `nullptr` in case of an allocation failure.

If an allocation has failed, a resumable function will use static member function `get_return_object_on_allocation_failure()` to construct the return value.

The following is an example of such specialization

```
namespace std {
    template <typename T, typename... Ts>
    struct coroutine_traits<kernel_mode_future<T>, Ts...> {
        template <typename... Us>
        static auto get_allocator(Us&&...) {return std::kernel_allocator<char>{}; }
        using promise_type = kernel_mode_resumable_promise<T>;

        static auto get_return_object_on_allocation_failure() { ... }
    };
}
```

Generalizing coroutine's promise `set_exception`

In exception-less environment, a requirement on `set_exception` member of coroutine promise needs to be relaxed to be able to take a value of an arbitrary error type `E` and not be limited to just the values of type `std::exception_ptr`. In exception-less environment, not only `std::exception_ptr` type may not be supported, but even if it were supported it is impossible to extract an error from it without relying on throw and catch mechanics.

Await expression: Unwrapping of an eventual value

As described earlier **await** *cast-expression* expands into an expression equivalent of:

```
{
    auto && __expr = cast-expression;
    if ( !await-ready-expr ) {
        await-suspend-expr;
        suspend-resume-point
    }
    cancel-check;
    return await-resume-expr;
}
```

Where `cancel-check` is expanded as

```
if ( promise-expr.cancellation_requested() ) goto end-label;
```

A straightforward implementation of `await_resume()` for getting an eventual value from the `future<T>` will call `.get()` that will either return the stored value or throw an exception. If unhandled, an exception

will be caught by a `catch(...)` handler of the resumable function and stored as an eventual result in a coroutine return object.

In the environments where exceptions are not allowed, implementation can probe for success or failure of the operation prior to resuming of the coroutine and use `<promise>.set_exception` to convey the failure to the promise. Coroutine promise, in this case, need to have `cancellation_requested()` to return true if an error is stored in the promise.

Here is how `await_suspend` may be defined for our hypothetical `kernel_future<T>` as follows:

```
template <typename Promise>
void kernel_future::await_suspend(std::coroutine_handle<Promise> p) {
    this->then([p](kernel_future<T> const& result) {
        if (result.has_error())
        {
            p.promise().set_exception(result.error());
        }
        p(); // resume the coroutine
    });
}
```

Appendix D demonstrates a complete implementation of adapters for `boost::future`, utilizing exception-less propagation technique described above.

Await expression: Failure to launch an asynchronous operation

If an `await_suspend` function failed to launch an asynchronous operation, it needs to prevent suspension of a resumable function at the await point. Normally, it would have thrown an exception and would have avoided suspend-resume-point. In the absence of exceptions, we can require that `await_suspend` must return false, if it failed to launch an operation and true otherwise. If false is returned from `await_suspend`, then coroutine will not be suspended and will continue execution. Failure can be indicate via `set_exception` mechanism as described in the previous section.

With all of the changes described in this section, `await expr` will be expanded into equivalent of:

```
{
    auto && __expr = cast-expression;
    if ( ! await-ready-expr && await-suspend-expr )
        suspend-resume-point
    }
    cancel-expression;
    return await-resume-expr;
}
```

With the extensions described above it is possible to utilize `await` and resumable functions in the environment where exceptions are banned or not supported.

Asynchronous cancellation

An attempt to cancel a coroutine that is currently suspended awaiting completion of an asynchronous I/O, can race with the resumption of a coroutine due to I/O completion. The coroutine model described in this paper can be extended to tackle asynchronous cancellation. Here is a sketch.

A coroutine promise can expose `set_cancel_routine(Fn)` function, where `Fn` is a function or a function object returning a value convertible to `bool`. A `set_cancel_routine` function should return `true` if `cancel_routine` is set and there is no cancellation in progress and `false` otherwise.

`await_suspend(std::coroutine_handle<Promise> rh)`, in addition to subscribing to get a completion of an asynchronous operation can use `rh.promise().set_cancel_routine(Fn)` to provide a callback that can attempt to cancel an asynchronous operation.

If a coroutine needs to be cancelled, it invokes a `cancel_routine` if one is currently associated with the coroutine promise. If `cancel_routine` returns `true`, it indicates that the operation in progress was successfully cancelled and the coroutine will not be resumed by the asynchronous operation. Thus, the execution of the coroutine is under full control of the caller. If `cancel_routine` returns `false`, it means that an asynchronous operation cannot be cancelled and coroutine may have already been resumed or will be resumed at some point in the future. Thus, coroutine resources cannot be released until pending asynchronous operation resumes the coroutine.

The following is an example of extending `sleep_for` awaiter from Appendix C to support asynchronous cancellation.

```
template <typename Promise>
bool await_suspend(std::coroutine_handle<Promise> resume_cb) {
    auto & promise = resume_cb.promise();
    if (promise.begin_suspend())
    {
        timer = CreateThreadpoolTimer(TimerCallback, resume_cb.to_address(), 0);
        if (timer)
        {
            promise.set_cancel_routine(timer, TimerCancel);
            SetThreadpoolTimer(timer, (PFILETIME)&duration, 0, 0);
            promise.done_suspend();
            return true;
        }
        else {
            promise.set_exception(
                std::system_error(std::system_category(), GetLastError()));
        }
    }
    promise.cancel_suspend();
    return false;
}
```

Where `begin_suspend`, `cancel_suspend` and `done_suspend` are used to help to solve races when cancellation is happening concurrently with invocation of `await_suspend`.

We do not propose this mechanism yet as we would like to gain more experience with developing libraries utilizing resumable functions described in this paper.

Stateful Allocators Support

Current proposal allows coroutines to be used with stackless and stackful allocators. To use a stateful allocator, coroutine traits's `get_allocator` need to select which parameters to a resumable function that carry an allocator to be used to allocate the coroutine state. Library designer can choose different

strategies, he/she could use `std::allocator_arg_t` tag argument followed by an allocator, or decide that allocator, if present, should be the first or the last parameter to a resumable function.

For example, using a generator coroutine from Appendix A and providing the following coroutine traits will enable stateful allocator use.

```
namespace std {
    template <typename R, typename Alloc, typename... Ts>
    struct coroutine_traits<generator<R>, allocator_tag_t, Alloc, Ts...> {
        template <typename... Us>
        static auto get_allocator(allocator_tag, Alloc a, Us&&...) {
            return a;
        }
        using promise_type = generator<R>::promise_type;
    };
}

template <typename Alloc>
generator<int> fib(allocator_tag_t, Alloc, int n)
{
    int a = 0;
    int b = 1;

    while (n-- > 0)
    {
        yield a;
        auto next = a + b;
        a = b;
        b = next;
    }
}

extern MyAlloc g_MyAlloc;

int main()
{
    for (auto v : fib(allocator_tag, g_MyAlloc, 35)) {
        std::cout << v << std::endl;
        if (v > 10)
            break;
    }
}
```

Override Selection of Coroutine Traits

In Urbana, several people asked for a way to select coroutine traits to use with the resumable function without altering the function signature. One possible syntax for `coroutine_traits` selection override could be as follows:

```
generator<int> fib() using(my_coroutine_traits)
{
    body
}
```

Proposed Standard Wording

No wording is provided at the moment.

Acknowledgments

Great thanks to Artur Laksberg, Chandler Carruth, Gabriel Dos Reis, Deon Brewis, James McNellis, Stephan T. Lavavej, Herb Sutter, Pablo Halpern, Robert Schumacher, Michael Wong, Niklas Gustafsson, Nick Maliwacki, Vladimir Petter, Slava Kuznetsov, Oliver Kowalke, Lawrence Crowl, Nat Goodspeed, Christopher Kohlhoff for your review and comments and Herb, Artur, Deon and Niklas for trailblazing, proposing and implementing resumable functions v1.

References

- [N3936] [Working Draft, Standard for Programming Language C++](#)
- [Revisiting Coroutines] [Moura, Ana Lúcia De and Ierusalimsky, Roberto. "Revisiting coroutines". ACM Trans. Program. Lang. Syst., Volume 31 Issue 2, February 2009, Article No. 6](#)
- [N3328] [Resumable Functions](#)
- [N3977] [Resumable Functions: wording](#)
- [N3985] [A proposal to add coroutines to the C++ standard library \(Revision 1\)](#)
- [Boost.Context] [Boost.Context Overview](#)
- [Boost.Coroutine] [Boost.Coroutine Overview](#)
- [SplitStacks] [Split Stacks in GCC](#)
- [JSF] [Join Strike Fighter C++ Coding Standards](#)
- [GoLang] <http://golang.org/doc/>
- [N3872] [A Primer on Scheduling Fork-Join Parallelism with Work Stealing](#)
- [Rx] <http://rx.codeplex.com/>
- [ReactiveX] <http://reactivex.io/>
- [RxAtNetflix] <http://techblog.netflix.com/2013/01/reactive-programming-at-netflix.html>
- [FRP] http://en.wikipedia.org/wiki/Functional_reactive_programming
- [PythonGenExprs] <http://legacy.python.org/dev/peps/pep-0289/>
- [BoostRangeAdapter] http://www.boost.org/doc/libs/1_49_0/libs/range/doc/html/range/reference/adaptors/introduction.html

Appendix A: An example of generator coroutine implementation

```
#include <resumable>
#include <iterator>

template <typename _Ty>
struct generator
{
    struct promise_type
    {
        enum class _StateT {_Active, _Cancelling, _Closed};

        _Ty const * _CurrentValue;
        _StateT _State = _StateT::_Active;

        promise_type& get_return_object() { return *this; }
        suspend_always initial_suspend() { return {}; }

        suspend_always final_suspend() {
            _State = _StateT::_Closed;
            return {};
        }

        bool cancellation_requested() const { return _State == _StateT::_Cancelling; }
        void set_result() {}

        suspend_always yield_value(_Ty const& _Value)
        {
            _CurrentValue = addressof(_Value);
            return {};
        }
    }; // struct generator::promise_type

    struct iterator : std::iterator<input_iterator_tag, _Ty>
    {
        coroutine_handle<promise_type> _Coro;

        iterator(nullptr_t): _Coro(nullptr) {}
        iterator(coroutine_handle<promise_type> _CoroArg) : _Coro(_CoroArg) {}

        iterator& operator++(){
            _Coro();
            if (_Coro.promise()._State == promise_type::_StateT::_Closed)
                _Coro = nullptr;
            return *this;
        }

        iterator operator++(int) = delete;
        // generator iterator's current_value is a reference to a temporary
        // on the coroutine frame. Implementing postincrement will require
        // storing a copy of the value in the iterator.

        bool operator==(iterator const& _Right) const { return _Coro == _Right._Coro; }
        bool operator!=(iterator const& _Right) const { return !(*this == _Right); }

        _Ty const& operator*() const {
            auto& _Prom = _Coro.promise();

```

```

        return *_Prom._CurrentValue;
    }
    _Ty const* operator->() const { return std::addressof(operator*()); }

}; // struct generator::iterator

iterator begin() {
    if (_Coro) {
        _Coro();
        if (_Coro.promise()._State == promise_type::_StateT::_Closed)
            return {nullptr};
    }
    return {_Coro};
}
iterator end() { return {nullptr}; }

explicit generator(promise_type& _Prom)
    : _Coro(coroutine_handle<promise_type>::from_promise(_STD addressof(_Prom)))
{}

generator() = default;
generator(generator const&) = delete;
generator& operator = (generator const&) = delete;

generator(generator && _Right): _Coro(_Right._Coro) { _Right._Coro = nullptr; }

generator& operator = (generator && _Right) {
    if (&_Right != this) {
        _Coro = _Right._Coro;
        _Right._Coro = nullptr;
    }
}

~generator() {
    if (_Coro) {
        auto& _Prom = _Coro.promise();
        if (_Prom._State == promise_type::_StateT::_Active) {
            // Note: on the cancel path, we resume the coroutine twice.
            // Once to resume at the current point and force cancellation.
            // Second, to move beyond the final_suspend point.
            _Prom._State = promise_type::_StateT::_Cancelling;
            _Coro();
        }
        _Coro();
    }
}

private:
    coroutine_handle<promise_type> _Coro = nullptr;
};

```

Appendix B: boost::future adapters

```
#include <resumable>
#define BOOST_THREAD_PROVIDES_FUTURE_CONTINUATION
#include <boost/thread/future.hpp>

namespace boost {
    template <class T>
    bool await_ready(unique_future<T> & t) { return t.is_ready();}

    template <class T, class Callback>
    void await_suspend(unique_future<T> & t, Callback cb)
    {
        t.then( [cb](auto&){ cb(); } );
    }

    template <class T>
    auto await_resume(unique_future<T> & t) { return t.get(); }
}

namespace std {

    template <class T, class... Whatever>
    struct coroutine_traits<boost::unique_future<T>, Whatever...> {
        struct promise_type
        {
            boost::promise<T> promise;

            auto get_return_object() { return promise.get_future(); }
            suspend_never initial_suspend() { return{}; }
            suspend_never final_suspend() { return{}; }

            template <class U = T, class = enable_if_t< is_void<U>::value >>
            void set_result() {
                promise.set_value();
            }

            template < class U, class U2 = T,
                        class = enable_if_t < !is_void<U2>::value >>
            void set_result(U&& value) {
                promise.set_result(std::forward<U>(value));
            }

            void set_exception(std::exception_ptr e){promise.set_exception(std::move(e));}

            bool cancellation_requested() { return false; }
        };
    };
}
```

Appendix C: Awaitable adapter over OS async facilities

```
#include <resumable>
#include <threadpoolapiset.h>

// usage: await sleep_for(100ms);
auto sleep_for(std::chrono::system_clock::duration duration) {
    class awaiter {
    public:
        static void TimerCallback(PTP_CALLBACK_INSTANCE, void* Context, PTP_TIMER) {
            std::coroutine_handle<>::from_address(Context)();
        }
        PTP_TIMER timer = nullptr;
        std::chrono::system_clock::duration duration;
    public:
        awaiter(std::chrono::system_clock::duration d) : duration(d){}
        bool await_ready() const { return duration.count() <= 0; }
        void await_suspend(std::coroutine_handle<> resume_cb) {
            int64_t relative_count = -duration.count();
            timer = CreateThreadpoolTimer(TimerCallback, resume_cb.to_address(), nullptr);
            if (timer == 0) throw std::system_error(GetLastError(), std::system_category());
            SetThreadpoolTimer(timer, (PFILETIME)&relative_count, 0, 0);
        }
        void await_resume() {}
        ~awaiter() {
            if (timer) CloseThreadpoolTimer(timer);
        }
    };
    return awaiter{ duration };
}
```

Appendix D: Exceptionless error propagation with boost::future

```
#include <boost/thread/future.hpp>

namespace boost {
    template <class T>
    bool await_ready(unique_future<T> & t) { return t.is_ready();}

    template <class T, class Promise>
    void await_suspend(
        unique_future<T> & t, std::coroutine_handle<Promise> rh)
    {
        t.then([=](auto& result){
            if(result.has_exception())
                rh.promise().set_exception(result.get_exception_ptr());
            rh();
        });
    }

    template <class T>
    auto await_resume(unique_future<T> & t) { return t.get(); }
}

namespace std {
    template <typename T, typename... anything>
    struct coroutine_traits<boost::unique_future<T>, anything...> {
        struct promise_type {
            boost::promise<T> promise;
            bool cancelling = false;
            auto get_return_object() { return promise.get_future(); }
            suspend_never initial_suspend() { return{}; }
            suspend_never final_suspend() { return{}; }
            template <class U> void set_result(U && value) {
                promise.set_value(std::forward<U>(value));
            }
            void set_exception(std::exception_ptr e) {
                promise.set_exception(std::move(e));
                cancelling = true;
            }
            bool cancellation_requested() { return cancelling; }
        };
    };
};
```