

Document No: N4402

Supersedes: N4286

Date: 2015-03-31

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

Resumable Functions (revision 4)

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 | 4 |
| Stackful Coroutine / Fiber / User-Mode thread | 4 |
| Split Stack / Linked Stack / Segmented Stack | 4 |
| Resumable Function | 4 |
| Suspend/Resume Point | 4 |
| Initial Suspend/Resume Point | 4 |
| Final Suspend/Resume Point | 4 |
| Discussion | 4 |
| Stackless vs. Stackful | 5 |
| Implementation Experience | 5 |
| Asynchronous I/O | 5 |
| Generator | 6 |
| Parent-stealing parallel computations | 7 |
| Go-like channels and goroutines | 7 |
| Reactive Streams | 8 |
| Resumable lambdas as generator expressions | 8 |
| Conceptual Model | 8 |
| Resumable Function | 8 |
| Resumable traits | 9 |

| | |
|--|----|
| Override Selection of Resumable traits | 10 |
| Allocation and parameter copy optimizations..... | 10 |
| auto / decltype(auto) return type..... | 10 |
| C++ does not need generator expressions... it already has them! | 11 |
| Resumable promise Requirements..... | 11 |
| Resumable handle..... | 12 |
| await operator | 14 |
| Evaluation of await expression | 14 |
| yield statement | 15 |
| Return statement..... | 16 |
| await-for statement | 17 |
| Resumable functions in environments where exceptions are unavailable or banned | 18 |
| Allocation failure..... | 18 |
| Generalizing coroutine's promise set_exception | 19 |
| Await expression: Unwrapping of an eventual value | 19 |
| Asynchronous cancellation | 20 |
| Stateful Allocators Support..... | 21 |
| Proposed Standard Wording..... | 22 |
| Acknowledgments..... | 22 |
| References | 22 |
| Appendix A: An example of generator coroutine implementation | 23 |
| Appendix B: boost::future adapters | 25 |
| Appendix C: Awaitable adapter over OS async facilities..... | 26 |
| Appendix D: Exceptionless error propagation with boost::future..... | 27 |

Revisions and History

This document supersedes N4286. Changes relative to N4286 include:

- renaming customization points back to resumable_traits and resumable_handle as they were in N4134;
- changing requirements on the return type of initial_suspend() and final_suspend() to be lexically convertible to bool;
- changing requirements on the return type of yield_value() to be either of a void type or a type lexically convertible to bool;

- making a `set_result` member function optional in a promise type; the absence of a `set_result` indicates that resumable function does not support an eventual return value and using `await` operator or `for-await` statement is not allowed in resumable functions with such promise;
- altered cancellation mechanism in resumable functions; instead of using `cancellation_requested()` member function in a promise to indicate that on the next resume resumable function need to be cancelled, an explicit member function `destroy()` is added to the `resumable_handle`. A `destroy()` member function can be invoked to force resumption of coroutine to go on the cancel path;
- added `resume()`, `destroy()`, and `done()` member functions to `resumable_handle`;
- moved proposed wording into a separate document N4302;
- removed trivial awaitables `suspend_always`, `suspend_never` and `suspend_if`; switching `yield_value`, `initial_suspend`, and `final_suspend` to return `bool` eliminated the need for those;

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.

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.

Suspend/Resume Point

A point at which execution of a resumable function can be suspended with a possibility to be resumed at a later time.

Initial Suspend/Resume Point

A suspend resume point that occurs prior to executing user-authored body of the resumable function.

Final Suspend/Resume Point

A suspend resume point that occurs after executing user-authored body of the resumable function, but, before resumable function state is destroyed.

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

¹ Allocation of the frame of a resumable function

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

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;
    }
}
```

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;
```

³ In constant space with only $\log n$ iterations

```
}
```

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 (auto v : input)
    {
        sum += v;
    }
    return sum;
}
```

As a producer:

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

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(auto 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.

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.

Resumable traits

Resumable 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::resumable_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 resumable_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 may define the primary template for resumable_traits as follows:

```
template <typename R, typename... Ts>
struct resumable_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 Resumable 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.

Override Selection of Resumable traits

In Urbana, a question was raised if there could be a way to alter selection of resumable traits without altering the function signature. One possible syntax for resumable_traits selection override could be as follows:

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

Where my_resumable_traits is a concrete type to be used in place of the resumable_traits specialization that would be normally chosen use the rules described earlier.

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 declspec(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>.

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(auto x:s) yield x*x; } ; // C++
```

Resumable 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 | A resumable promise type |
| p | A value of type P |
| e | A value of std::exception_ptr type |
| h | A value of type std::resumable_handle<P> |
| v | An <i>expression</i> or <i>braced-init-list</i> |

Table 2: Coroutine Promise Requirements

| Expression | Note |
|-----------------------|---|
| P{} | Construct an object of type P. |
| p.get_return_object() | get_return_object is invoked by the resumable function to construct the return object prior to reaching the first suspend-resume point, a return statement, or flowing off the end of the function. |
| p.set_result(v) | If present, an enclosing resumable function supports an eventual value of a type that v can be converted to. set_result is invoked by a |

| | |
|--------------------|--|
| | resumable function when a return statement with an <i>expression</i> or a <i>braced-init-list</i> is encountered in a resumable function. If a promise type does not satisfy this requirement, the presence of a return statement with an expression or a <i>braced-init-list</i> statement in the body results in a compile time error. |
| p.set_result() | If present, an enclosing resumable function supports an eventual value of type void. set_result is invoked by a resumable function when a return statement without an <i>expression</i> nor a <i>braced-init-list</i> is encountered in a resumable function or control flows to the end of the function. A promise type must contain at most one declaration of set_result. |
| p.set_exception(e) | set_exception is invoked by a resumable function when an unhandled exception occurs within a body of the resumable function. If promise does not provide set_exception, unhandled exceptions will propagate from a resumable functions normally. |
| p.yield_value(v) | Must be present for the enclosing resumable function to support yield statement. |
| initial_suspend() | if p.initial_suspend() evaluates to true, resumable function will suspend at initial suspend point |
| final_suspend() | if p.final_suspend() evaluates to true, resumable function will suspend at final suspend point |

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.

Resumable handle

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 member functions of a resumable handle referring to that function.

Synopsis:

```
template <typename Promise = void> class resumable_handle;

template<> struct resumable_handle<void>
{
    // construct/reset
    resumable_handle()noexcept;
    resumable_handle(std::nullptr_t) noexcept;
    resumable_handle& operator=(nullptr_t) noexcept;

    // export/import
    static resumable_handle from_address(void* addr) noexcept;
    void* to_address() const noexcept;

    // capacity
    explicit operator bool() const noexcept;
```

```

    // resumption
    void operator>()() const;
    void resume() const;
    void destroy() const;

    // completion check
    bool done() const noexcept;
};

template <typename Promise>
struct resumable_handle : resumable_handle<>
{
    // construct/reset
    using resumable_handle<>::resumable_handle;
    resumable_handle& operator=(nullptr_t) noexcept;

    // export/import
    static resumable_handle from_promise(Promise*) noexcept;
    Promise& promise() noexcept;
    Promise const& promise() const noexcept;
};

```

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

Resumable handle has two forms. One that provides an ability to resume evaluation of a resumable function and another, which additionally allows access to the promise of the associated resumable function.

In this revision, we added three new member functions: resume, destroy, and done.

- resume() member function has the same effect as operator(), namely, it resumes the function at the current suspend point.
- destroy() member function that “cancel”s the resumable function, by destroying objects with automatic storage duration that are in scope at the suspend point in the reverse order of the construction. If resumable function required dynamic allocation for the objects with automatic storage duration, the memory is freed.

This change eliminates the need for cancellation_requested() member of resumable promise, simplifies authoring new coroutine types and allows the optimizer to reason about the lifetime of the resumable function.

- done() member function indicates whether the coroutine is suspended at final suspend/resume point and eliminates the need for separate state tracking in the promise of the coroutine.

Bikeshed: resumption_handle, resumption_object, resumable_ptr, basic_resumable_handle instead of resumable_handle<void>, from_raw_address, to_raw_address, from_pvoid, to_pvoid.

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

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
    }
    return await-resume-expr;
}
```

if the type of await-suspend-expr is void, otherwise it is equivalent to

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

Where __expr is a variable defined for exposition only, and _ExprT is the type of the *cast-expression*, and _ResumableHandle is an object of the resumable_handle type associated with the enclosed function and the rest are determined as follows:

| | |
|---|---|
| await-ready-expr await-suspend-expr await-resume-expr | <p>— if _ExprT is a class type, the unqualified-ids await_ready, await_suspend and await_resume are looked up in the scope of class _ExprT as if by class member access lookup (N3936/3.4.5 [basic.lookup.classref]), and if it finds at least one declaration, then await_ready-expr, await_suspend-expr and await_resume-expr are __expr.await_ready(), __expr.await_suspend(resumption-function-object) and __expr.await_resume(), respectively;</p> <p>— otherwise, await_ready-expr, await_suspend-expr and await_resume-expr are await_ready(__expr), await_suspend(__expr, resumption-function-object) and await_resume(__expr), respectively, where await_ready, await_suspend and await_resume are looked up in the</p> |
|---|---|

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.

| | |
|--|--|
| | <p>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> |
|--|--|

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.

Design Note: Motivation for the second form of `await_suspend` (with non void return type) is to allow efficient handling of cases where asynchronous operation is launched in `await_suspend` and there is an indication from the OS that async operation was not launched as it has already completed synchronously, in this case, there is no need to suspend, but, we don't know that until we invoke OS API to launch an asynchronous operation. In earlier papers this form was also motivated by using resumable functions in exception-less environment, that use case is no longer needed due to simplified cancellation model proposed in this paper.

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

yield statement

A yield statement of the form

yield V;

where V is either an *expression* or a *braced-init-list* is equivalent to:

```
_Pr.yield_value(V);
suspend-resume-point
```

If a `_Pr.yield_value(V)` expression is of type **void**, otherwise it is equivalent to:

```

    if (_Pr.yield_value(V)) {
        suspend-resume-point
    }

```

Where `_Pr` is a promise object of the enclosing resumable function.

A **yield** statement may only appear if `yield_value` member function is defined in the promise type of the enclosing resumable function.

A promise object may have more than one overload of `yield_value`. Consider this example:

```

recursive_generator<int> flatten(node* n)
{
    if (n == nullptr)
        return;

    yield flatten(n->left);
    yield n->value;
    yield flatten(n->right);
}

```

In the example above, the promise for `flatten` function should contain overloads that can accept a value of type `int` and a value of type `recursive_generator<int>`. In the former case, yielding a value is unconditional. In the latter case, the nested generator may produce an empty sequence of values and thus suspension at the yield point shall not happen and corresponding `yield_value` contextually converted to `bool` must evaluate to false.

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.

In the experimental implementation in Microsoft Visual C++ we chose to use the following rule: if `yield` token is followed by `'('`, `yield` is an identifier, otherwise, it is a keyword.

Return statement

Let `_Pr` be a promise object of a resumable function, `v` be an *expression* or *braced-init-list*, and *end-label* is a label just before final-suspend-point of the enclosing resumable function.

If an expression `_Pr.set_result(v)` is valid, then, the resumable function supports an eventual value of some non-void type and a return statement in a resumable function in a form **return** `v`; is equivalent to:


```
{ _Pr.set_result(v); goto end-label; }
```

If an expression `_Pr.set_result()` is valid, then, the resumable function supports an eventual value of type `void` and a return statement in a resumable function in a form **return**; is equivalent to:

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

If `_Pr` has no overloads of `set_result`, then, the resumable function does not support an eventual value of any type and the only allowed form of the return statement is **return**; which is equivalent to:

```
goto end-label;
```

If a 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 operator and/or for-await statements are only allowed in resumable functions supporting an eventual value of some type.

Resumable function should have at most one overload of `set_result`.

Design note: In earlier revisions of this paper, `set_result` was a mandatory member of the promise object. We made it optional in this revision and tied to a notion of eventual return type that makes intuitive sense for asynchronous resumable functions (the ones with await operators and for-await statements) and less sense in synchronous generators.

The motivation for the change was driven by customer mistakes where they would accidentally use an await in a body of the generator resumable function that will result in suspension without assigning the next value to a generator, resulting in a crash or reading a garbage value.

With this change, generator's promise won't define `set_result` and thus using await or for-await in the body of the generator will be caught at compile time.

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.

Resumable functions in environments where exceptions are unavailable or 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, `resumable_traits` specialization can declare an optional static member function `get_return_object_on_allocation_failure`.

If a `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 a 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 resumable_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
    }
    return await-resume-expr;
}
```

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.

An `await_suspend` may be defined for our hypothetical `kernel_future<T>` as follows:

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

```

    }
    else
        p.resume();
    });
}

```

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

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::resumable_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>
void await_suspend(std::resumable_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();
        }
        else {
            promise.set_exception(
                std::system_error(std::system_category(), GetLastError()));
        }
        return;
    }
    promise.cancel_suspend();
}

```

```
}
```

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, Resumable 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 Resumable traits will enable stateful allocator use.

```
namespace std {
    template <typename R, typename Alloc, typename... Ts>
    struct resumable_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;
    }
}
```

Proposed Standard Wording

Wording is provided in a separate document N4403.

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, Shahms King, Slava Kuznetsov, Tongari J, 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
    {
        _Ty const * _CurrentValue;

        promise_type& get_return_object() { return *this; }
        bool initial_suspend() { return true; }
        bool final_suspend() { return true; }
        void yield_value(_Ty const& _Value) { _CurrentValue = addressof(_Value); }
    }; // struct generator::promise_type

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

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

        iterator& operator++(){
            _Coro.resume();
            if (_Coro.done())
                _Coro = nullptr;
            return *this;
        }

        iterator operator++(int) = delete; // ban postincrement to keep iterator lean by
            // avoiding storing the current value in the iterator itself

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

        _Ty const& operator*() const { return *_Coro.promise()._CurrentValue; }
        _Ty const* operator->() const { return std::addressof(operator*()); }
    }; // struct generator::iterator

    iterator begin() {
        _Coro.resume();
        if (_Coro.done())
            return {nullptr};
        return {_Coro};
    }
    iterator end() { return {nullptr}; }

    explicit generator(promise_type& _Prom)
        : _Coro(resumable_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;
    }
    return *this;
}
~generator() { if (_Coro) _Coro.destroy(); }
private:
    resumable_handle<promise_type> _Coro;
};

```


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 resumable_traits<boost::unique_future<T>, Whatever...> {
        struct promise_type
        {
            boost::promise<T> promise;

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

            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));}

        };
    };
}
```

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::resumable_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::resumable_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) {
        if(result.has_exception()) {
            return false; // returning true, would send us to await_ready()
                          // that will call .get() that will result
                          // in exception throwing. Returning false, passes
                          // control to await_suspend that will handle
                          // the error propagation without exceptions
        }
        return t.is_ready();
    }

    template <class T, class Promise>
    void await_suspend(
        unique_future<T> & t, std::resumable_handle<Promise> rh)
    {
        if(t.has_exception()) {
            rh.promise().set_exception(t.get_exception_ptr());
            rh.destroy(); // destroys the coroutine
        }
        else {
            t.then([=](auto& result){
                if(result.has_exception()) {
                    rh.promise().set_exception(result.get_exception_ptr());
                    rh.destroy(); // destroys the coroutine
                }
                else
                    rh.resume();
            });
        }
    }

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

namespace std {
    template <typename T, typename... anything>
    struct resumable_traits<boost::unique_future<T>, anything...> {
        struct promise_type {
            boost::promise<T> promise;
            auto get_return_object() { return promise.get_future(); }
            bool initial_suspend() { return false; }
            bool final_suspend() { return false; }
            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));
            }
        };
    };
};
```