

**Document Number:** N3747  
**Date:** 2013-09-01  
**Reply To** Christopher Kohlhoff <chris@kohlhoff.com>

---

# A UNIVERSAL MODEL FOR ASYNCHRONOUS OPERATIONS

## 1 INTRODUCTION

Networking has a long history of programs being written using event-driven and asynchronous programming designs. The use of asynchronous operations with continuations, in particular, offers a good model for abstraction and composition. Asynchronous operations may be chained, with each continuation initiating the next operation. The composed operations may then be abstracted away behind a single, higher level asynchronous operation with its own continuation.

N3634 proposes an evolution of `std::future` with the addition of a `then()` member function, amongst other things. This function allows one to attach a continuation to a future, and is intended for use with asynchronous operations. With these extensions, `std::future` is proposed as a standard model for representing asynchronous operations.

N3650 builds on these extensions to `std::future` with a new language facility, resumable functions. The new keywords, `async` and `await`, are intended to ease composition of asynchronous operations by enabling the use of imperative flow control primitives.

In this paper, we will first examine how futures can be a poor choice as a fundamental building block of asynchronous operations in C++. The extended `std::future` brings inherent costs that penalise programs, particularly in domains where C++ use is prevalent due to its low overheads. An asynchronous model based on a pure callback approach, on the other hand, allows efficient composition of asynchronous operations.

However, cognizant that some C++ programmers may have valid reasons for preferring a futures-based approach, this paper introduces a universal model for asynchronous operations. This model supports both lightweight callbacks and futures, allowing the application programmer to select an approach based on appropriate trade-offs.

Finally, we will see how this universal model can be leveraged to support other models of composition. This paper presents implementation experience of the universal model, which includes several pure library implementations of resumable functions, or coroutines. Programmers have the opportunity to express asynchronous control flow in an imperative manner, without requiring the addition of new keywords to the language.

## 2 CALLBACKS VERSUS FUTURES

This paper uses the terms *callbacks* and *futures* as shorthand for two asynchronous models.

These two models have several concepts in common:

- An *initiating function* that starts a given asynchronous operation. The arguments to this function, including an implicit `this` pointer, supply the information necessary to perform the asynchronous operation.
- A *continuation* that specifies code to be executed once the operation completes.

The *callbacks* model<sup>1</sup> refers to a design where the continuation is passed, in the form of a function object, as an argument to the initiating function:

```
socket.async_receive(args, continuation);
```

In this example, `async_receive()` is the initiating function. When the asynchronous receive operation completes, the result is passed as one or more arguments to the callback object `continuation`. In Boost.Asio, these continuations are called *handlers*.

With the *futures* model, the initiating function returns a future object. The caller of the initiating function may then attach a continuation to this future:

```
socket.async_receive(args).then(continuation);
```

Or, more explicitly:

```
std::future<size_t> fut = socket.async_receive(args);  
...  
fut.then(continuation);
```

Alternatively, the caller may choose to do something else with the future, such as perform a blocking wait or hand it off to another part of the program. This separation of the asynchronous operation initiation from the attachment of the continuation can sometimes be a benefit of the futures model. However, as we will see below, this separation is also the source of the runtime costs associated with the model.

To maximise the usefulness of a C++ asynchronous model, it is highly desirable that it support callbacks. Reasons for preferring a callbacks model include:

- Better performance and lower abstraction penalty.
- A fundamental building block. Futures, resumable functions, coroutines, etc. can be efficiently implemented in terms of callbacks.
- Not tied to threading facilities. It is possible to implement efficient callback-based network programs on platforms that have no thread support.

This paper aims to demonstrate that supporting callbacks need not be mutually exclusive to providing support for futures.

---

<sup>1</sup> An implementation of the callbacks model in C++ can be found in the Boost.Asio library, in Boost versions up to and including 1.53.

### 3 PERFORMANCE MATTERS

Long-time network programmers may have encountered the following misconception: that performance is not a primary concern, due to the high latencies involved in network I/O.

Latency may be an acceptable justification for a high abstraction penalty in the context of HTTP clients on desktop operating systems. However, higher overheads mean lower throughput. For programs, such as servers, that handle network I/O events from multiple sources, the latency to the peer is often irrelevant; the program needs to be ready to handle the next event immediately. Furthermore, the deleterious effects of queuing and congestion are felt well before a system reaches 100% utilisation.

And, while it is true that typical Internet latencies are high, often measured in tens or hundreds of milliseconds, high latency is not an inherent attribute of network I/O. By way of illustration, consider the following two data points:

⇒ Transmit a 64-byte UDP packet from a user-space application on one host to a user-space application on another host (i.e. `RTT/2`), across a 10GbE network.

**2 microseconds**

⇒ On the same hardware and operating system, wake up and switch to a thread that is blocked on a `std::future` object.

**3 microseconds**

There are many real world use cases where C++ is used because it allows for high-level abstractions with a low abstraction penalty. For example, the author is familiar with systems in the financial markets domain where performance differences measured in microseconds have a significant impact on an application's efficacy.

It is for these use cases that the choice of asynchronous model matters most. If C++ were to adopt a restricted asynchronous model based only on futures, potential C++ standard library components such as networking would have their usefulness limited. To meet their application's performance requirements, programmers of these systems would have to step outside the standard library. Put another way, C++ and its standard library would have less to differentiate it from other, higher-level languages.

### 4 ANATOMY OF AN ASYNCHRONOUS OPERATION

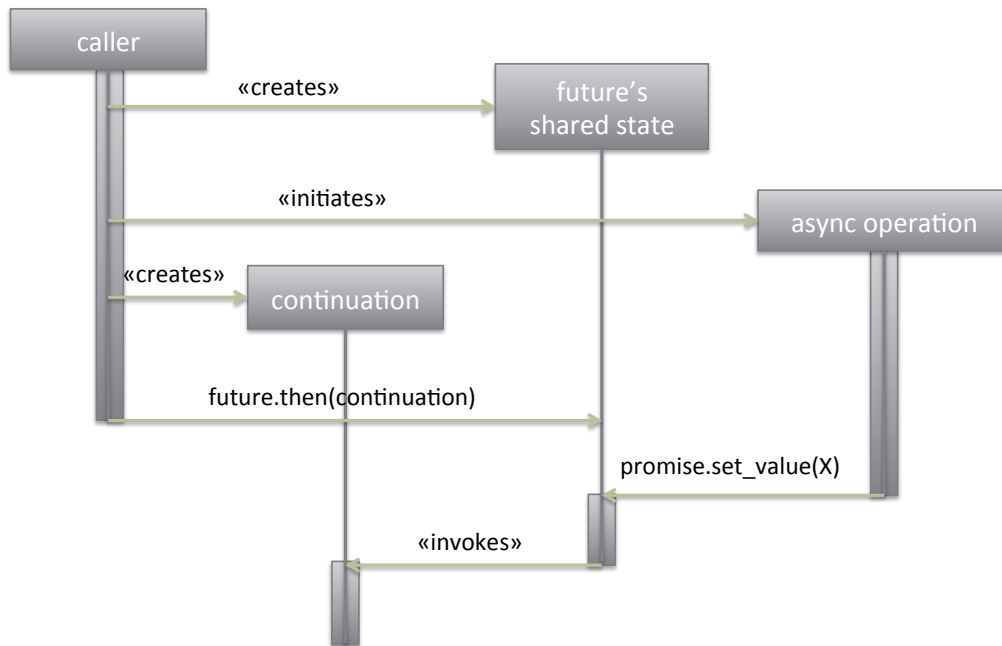
To see the inherent costs of a futures model, let us take code of the form:

```
socket.async_receive(args).then(continuation);
```

and consider the underlying sequence of events<sup>2</sup>:

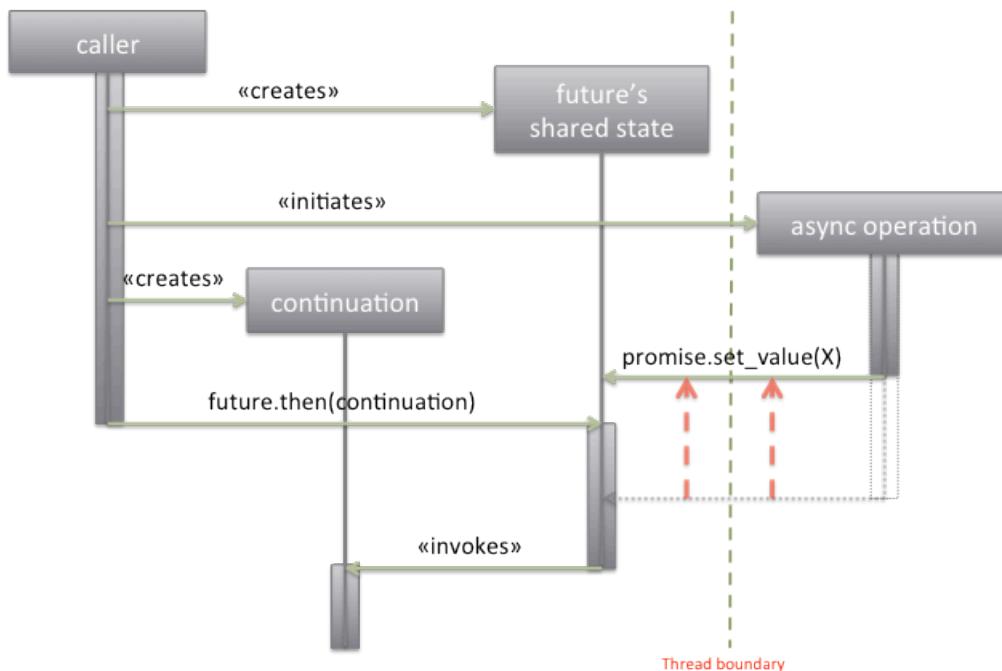
---

<sup>2</sup> In reality, the implementation is more complicated than presented, as the `then()` member function itself returns an additional `std::future` object with another shared state. It is likely that this second shared state is where the continuation is ultimately stored.



The initiating function creates the future's shared state and launches the asynchronous operation. The caller then attaches the continuation to the future. Some time later, the asynchronous operation completes and the continuation is invoked.

However, after the asynchronous operation is initiated, it is logically executing in its own thread of control. It is executing in parallel to the caller, and so it is possible for the operation to complete before the continuation is attached, as shown in the following sequence of events:

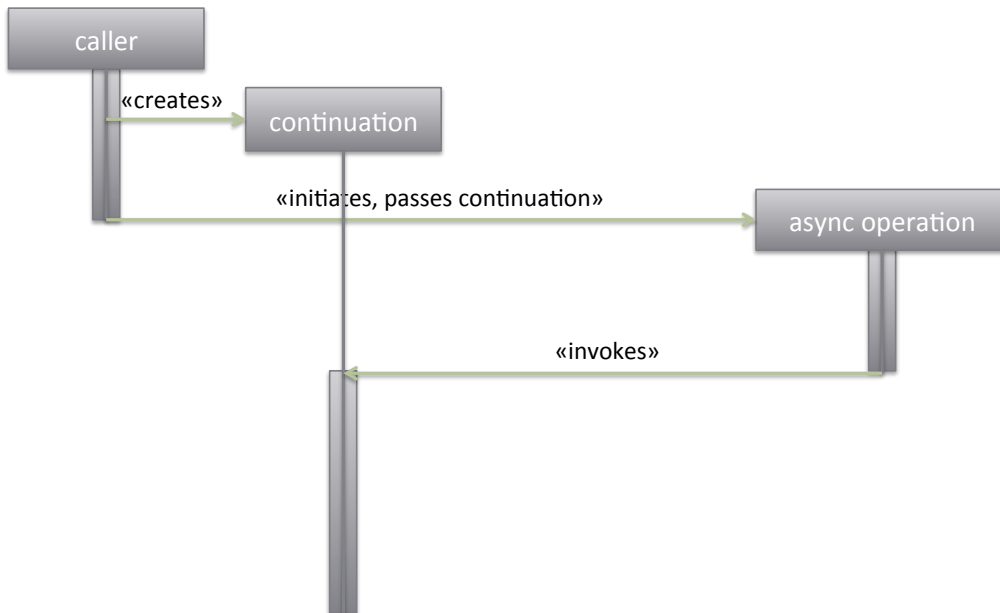


As a consequence, the shared state object has a non-deterministic lifetime, and requires some form of synchronisation to coordinate the attachment and invocation of the continuation.

In contrast, when the callbacks model has code of the form:

```
socket.async_receive(args, continuation);
```

we see the following, simpler sequence of events:



The initiating function accepts the continuation object and launches the asynchronous operation. The caller's flow of control ceases at this point<sup>3</sup>, and the asynchronous operation is not executing in parallel with it. Unlike the futures model, there are no shared objects with non-deterministic lifetime, and no additional synchronisation is required.

## 5 COMPOSITION OF OPERATIONS

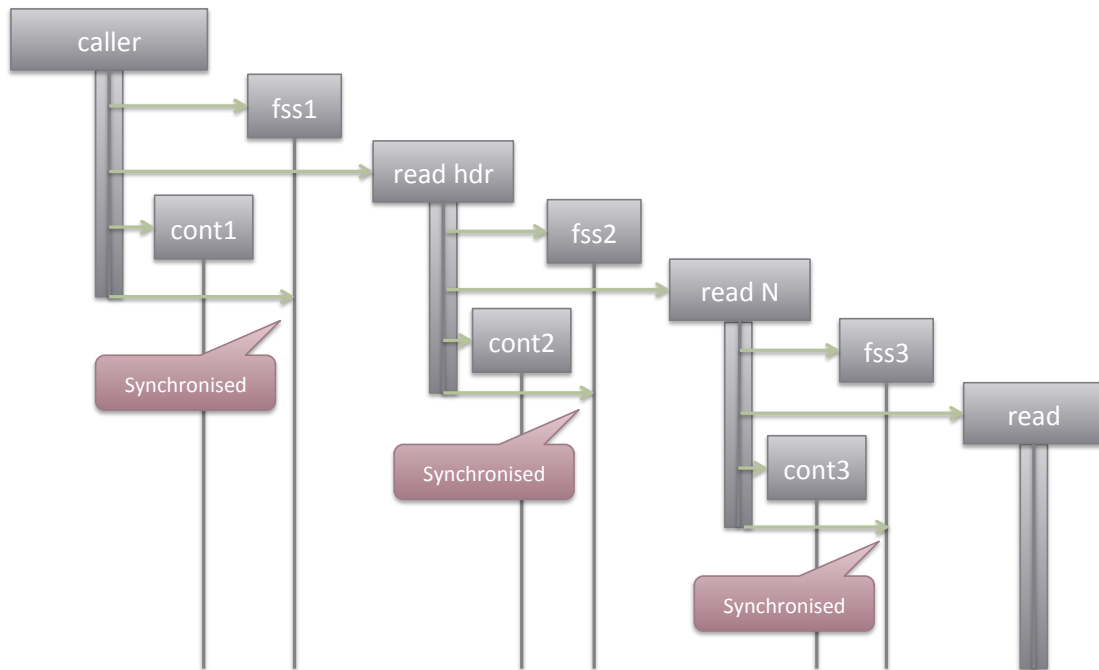
Let us now consider what happens when asynchronous operations are composed. A typical scenario involves a network protocol with a fixed-length header and variable-length body. For this example, the tree of operations might be as follows:

- ⇒ Read message
  - ⇒ Read header
    - ⇒ Read N bytes
      - ⇒ Read data off socket
  - ⇒ Read body
    - ⇒ Read N bytes
      - ⇒ Read data off socket

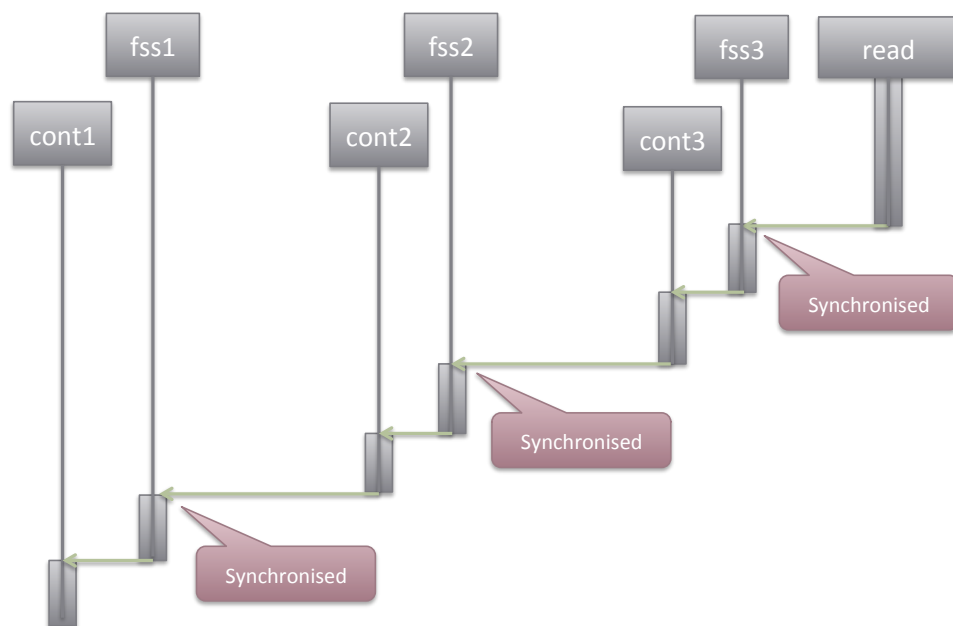
Each operation represents a level of abstraction, and has its own set of post-conditions that must be satisfied before its continuation can be invoked. For example, the "Read N bytes" operation exists to manage the problem of partial reads (where a socket returns fewer bytes than requested), and cannot call its continuation until the requested number of bytes is read or an error occurs.

With futures, as we go down the tree we see a sequence of events similar to the following:

<sup>3</sup> Technically, the caller's lifetime is not required to end at this time. It can continue to perform other computations or launch additional asynchronous operations. The important point is that it is not required to continue in parallel to the asynchronous operation in order to attach a continuation.

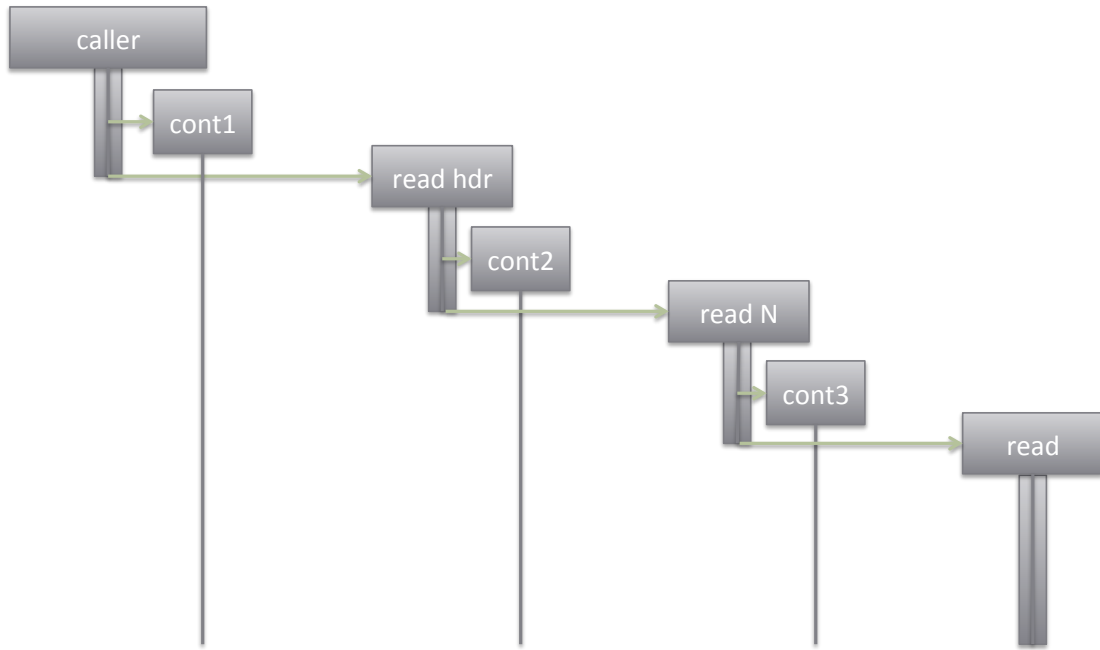


When the operations complete, and assuming each post-condition is immediately satisfied, the sequence of events is:

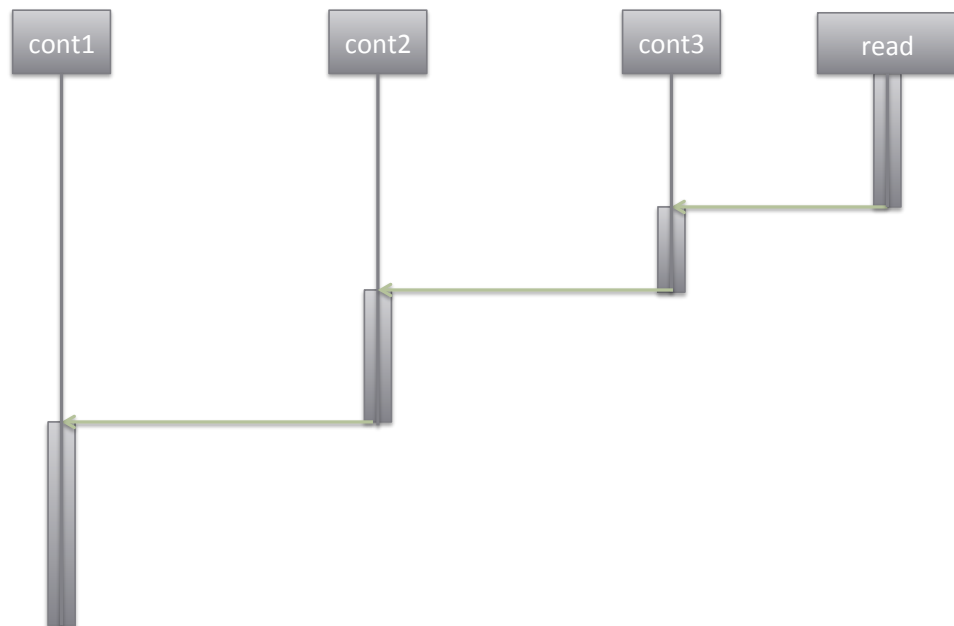


As you can see, each layer of the abstraction adds synchronisation overhead. On the systems available to the author, the cost of synchronisation mechanisms such as a `mutex` (when uncontended) or sequentially consistent atomic variable is some 10-15 nanoseconds per use.

When composing operations with callbacks, these costs are not incurred on the way down the tree of operations:



Nor are they present on the way back up:



In addition to avoiding synchronisation costs at each layer, when written as templates the compiler is given the opportunity to further optimise, such as inlining continuations. This means that it is possible to create layers of abstraction that have little or no runtime abstraction penalty.

## 6 A PLETHORA OF FUTURES

Let us limit ourselves, for a moment, to considering only a futures model. All asynchronous operations would take the following form:

```
class socket {
    ...
    std::future<size_t> async_receive(buffer_type buffer);
    ...
};
```

This approach raises the following questions:

- How do we specify a custom allocator?
- How do we support alternate future types? For example:
  - A future that waits / does not wait on destruction.
  - A future with no blocking operations at all.
  - A future for use only in single-threaded programs.
  - A future that works with 3<sup>rd</sup> party coroutines.
  - A future that works with a 3<sup>rd</sup> party executor.

In acknowledging that there are use cases where there are valid reasons for preferring a futures model to a callbacks model, the question is: can we have a single model that supports callbacks *and* futures (in all their incarnations)?

## 7 A UNIVERSAL MODEL

In Boost 1.54, the Boost.Asio library introduced a universal asynchronous model that supports callbacks, futures, and resumable functions or coroutines. The model is also user extensible to allow the inclusion of other facilities. Before we look at how the model works, in this section we will review some of the ways in which the model can be used.

### 7.1 CALLBACKS

Callbacks continue to work as before. As they are simply function objects, they can be specified using function pointers:

```
void handle_receive(error_code ec, size_t n) { ... }
...
socket.async_receive(buffer, handle_receive);
```

or lambdas:

```
socket.async_receive(buffer,
    [](error_code ec, size_t n)
    {
        ...
    });
```

or with function object binders:



```
socket.async_receive(buffer,  
    std::bind(handle_receive, _1, _2));
```

With some suitable macro magic<sup>4</sup>, we can even have callbacks that are implemented as “stackless” coroutines:

```
struct handler : coroutine  
{  
    ...  
    void operator()(error_code ec, size_t n)  
    {  
        reenter (this)  
        {  
            ...  
            while (!ec)  
            {  
                yield socket.async_receive(buffer, *this);  
                ...  
            }  
            ...  
        }  
    }  
    ...  
};
```

Here we can make use of imperative control flow structures to implement complex asynchronous logic in a synchronous manner. Runtime overhead is minimal: the coroutine state is stored in an integer, and re-entering a coroutine is equivalent to resuming a switch-based state machine.

## 7.2 FUTURES

By passing a special value `use_future` (similar in concept to how the global object `std::nothrow` is used to tag overloaded functions), initiating functions return a `std::future` that can be used to wait for the asynchronous operation to complete:

```
try  
{  
    future<size_t> n =  
        socket.async_receive(  
            buffer, use_future);  
    // Use n.get() to obtain result  
}  
catch (exception& e)  
{  
    ...  
}
```

The correct return type is automatically deduced based on the asynchronous operation being used. For example, this code calls `async_receive`, and the result is a `future<size_t>` to represent the number of bytes transferred. In section 8 we will see how this type is obtained.

---

<sup>4</sup> Implemented in terms of a `switch` statement using a technique similar to Duff’s Device.

If the asynchronous operation fails, the `error_code` is converted into a `system_error` exception and passed back to the caller through the future.

The `use_future` special value also allows a custom allocator to be specified:

```
try
{
    future<size_t> n =
        socket.async_receive(
            buffer, use_future[my_allocator]);
    // Use n.get() to obtain result
}
catch (exception& e)
{
    ...
}
```

The `use_future` object may likewise be passed to asynchronous operations that are themselves compositions. If these compositions are built using callbacks, the intermediate operations and their continuations are executed efficiently as in the callbacks model. Only at the final step is the future made ready with the result. However, should any intermediate step result in an exception, that exception is caught and stored on the future, where it will be re-thrown when the caller performs `get()`.

### 7.3 COROUTINES / RESUMABLE FUNCTIONS

Support for “stackful” coroutines has been implemented on top of the `Boost.Coroutine` and `Boost.Context` libraries. This is a pure library solution of resumable functions that does not require the addition of any new keywords.

A `yield_context` object is used to represent the current coroutine. By passing this object to an initiating function, we indicate that the caller should be suspended until the operation is complete:

```
void receive_message(yield_context yield)
{
    try
    {
        size_t n = socket.async_receive(buffer, yield);
        ...
    }
    catch (exception& e)
    {
        ...
    }
}
```

The return type of the initiating function is deduced based on the operation being called. If the operation fails, the `error_code` is converted into a `system_error` exception and passed back to the caller through the future. In many use cases, an error is not exceptional, and it is preferable that it be handled using other control flow mechanisms. With these coroutines, the error can be captured into a local variable:

```

void receive_message(yield_context yield)
{
    ...
    error_code ec;
    size_t n = socket.async_receive(buffer, yield[ec]);
    if (ec) ...
}

```

As each coroutine has its own stack, local variables and complex control flow structures are available, exactly as they would be in a synchronous implementation of the algorithm:

```

void do_echo(yield_context yield)
{
    try
    {
        char data[128];
        for (;;)
        {
            size_t n = socket.async_read_some(buffer(data), yield);
            async_write(socket, buffer(data, n), yield);
        }
    }
    catch (std::exception& e)
    {
        // ...
    }
}

```

A `yield_context` object may be passed to composed operations that are built only using callbacks. The coroutine functions themselves also compose easily through direct function calls. These functions share a stack, and the bottommost function suspends the coroutine until an asynchronous operation completes. Unlike futures, returning a result from a lower abstraction layer has minimal cost; it is the same as returning a result from a normal function.

Finally, as a comparison, here is an example of an `async/await`-based resumable function shown side-by-side with its equivalent using Boost.Asio's coroutines:

*Example using Microsoft's PPL<sup>5</sup>*

```

task<string>
read(string file, string suffix)
    __async {
    istream fi = __await open(file);
    string ret, chunk;
    while((chunk = __await fi.read()).size())
        ret += chunk + suffix;
    return ret;
}

```

*Equivalent using library-based coroutines*

```

string
read(string file, string suffix,
    yield_context yield) {
    istream fi = open(file, yield);
    string ret, chunk;
    while((chunk = fi.read(yield)).size())
        ret += chunk + suffix;
    return ret;
}

```

## 7.4 PROPOSED BOOST.FIBER LIBRARY

Boost.Fiber<sup>6</sup>, a library that is being developed and proposed for inclusion in Boost, provides a framework for cooperatively scheduled threads. As with "stackful" coroutines, each fiber has

<sup>5</sup> Copied from <http://video.ch9.ms/sessions/build/2013/2-306.pptx>

its own stack and is able to suspend its execution state. The Boost.Fiber library supplies many concurrency primitives that mirror those in the standard library, including mutexes, condition variables and futures.

The proposed Boost.Fiber library has been extended to support the universal asynchronous model. This has been achieved without requiring any Boost.Fiber-specific extensions to Boost.Asio.

Firstly, we can use Boost.Fiber's future class in the same way as `std::future`, except that waiting on the future suspends the current fiber:

```
try
{
    boost::fibers::future<size_t> n =
        socket.async_receive(buffer,
            boost::fibers::asio::use_future);
    // Use n.get() to obtain result
}
catch (exception& e)
{
    ...
}
```

Secondly, we can suspend the current fiber automatically when performing an asynchronous operation, in a similar fashion to the integration with Boost.Coroutine shown above:

```
try
{
    size_t n = socket.async_receive(buffer,
        boost::fibers::asio::yield);
    ...
}
catch (exception& e)
{
    ...
}
```

In both cases, the `use_future` and `yield` names refer to special values, similar to `std::nothrow`. When passed, the appropriate return type is deduced based on the asynchronous operation that is being called.

## 8 HOW THE UNIVERSAL MODEL WORKS

To understand the universal asynchronous model, let us first consider the callbacks model from Boost.Asio (used in Boost version 1.53 and earlier). We will then examine the incremental changes that have been applied to create the universal model.

In a callbacks model, a typical initiating function will look something like this:

---

<sup>6</sup> <https://github.com/olk/boost-fiber>

```
template <class Buffers, class Handler>
void socket::async_receive(Buffers b, Handler&& handler)
{
    ...
}
```

To convert to the universal model, we need to determine the return type and how to obtain the return value:

```
template <class Buffers, class Handler>
???? socket::async_receive(Buffers b, Handler&& handler)
{
    ...
    return ???;
}
```

This is achieved by introducing two customisation points into the implementation.

### 8.1 CUSTOMISATION POINT 1: CONVERT TO THE “REAL” HANDLER TYPE

In the callbacks model, the type `Handler` was the function object type to be invoked on completion of the operation. In the universal model, it may actually be a placeholder type, such as `yield_context` or the type of the `use_future` object.

Therefore it is first necessary to determine the real type of the handler. This is achieved by the following type trait:

```
template <typename Handler, typename Signature>
struct handler_type {
    typedef ... type;
};
```

The `Signature` template parameter is based on the callback arguments for the given asynchronous operation. For a socket receive operation the `Signature` is `void(error_code, size_t)`, and the real handler type may be deduced as follows:

```
typename handler_type<Handler, void(error_code, size_t)>::type
```

The real handler type must support construction from the placeholder type, as the initiating function will attempt to construct a real handler as follows:

```
typename handler_type<
    Handler, void(error_code, size_t)>::type
    real_handler(std::forward<Handler>(handler));
```

The `handler_type` template would be specialised for any type that must participate in the universal asynchronous model. For example, when we write:

```
auto fut = socket.async_receive(buffers, use_future);
```

the initiating function performs the equivalent of:

```
typename handler_type<
    use_future_t, void(error_code, size_t)>::type
    real_handler(handler);
```

which yields a `real_handler` object of type `promise_handler<size_t>`. The `promise_handler<>` template is an implementation detail of Boost.Asio, and it simply sets a `std::promise<>` object's value when an asynchronous operation completes.

## 8.2 CUSTOMISATION POINT 2: CREATE THE INITIATING FUNCTION'S RESULT

With the callbacks model, initiating functions always have a `void` return type. In the universal model, the return type must be deduced and the return value determined. This is performed through the `async_result` type:

```
template <typename Handler>
class async_result {
public:
    // The return type of the initiating function.
    typedef ... type;

    // Construct an async result from a given handler.
    explicit async_result(Handler& handler) { ... }

    // Obtain initiating function's return type.
    type get() { return ...; }
};
```

The `async_result` template is specialised for real handler types, and acts as the link between the handler (i.e. the continuation) and the initiating function's return value. For example, to support `std::future` support the template is specialised for the `promise_handler<>` template:

```
template <typename T>
class async_result<promise_handler<T>> {
public:
    // The return type of the initiating function.
    typedef future<T> type;

    // Construct an async result from a given handler.
    explicit async_result(promise_handler<T>& h) { f_ = h.p_.get_future(); }

    // Obtain initiating function's return value.
    type get() { return std::move(f_); }

private:
    future<T> f_;
};
```

## 8.3 PUTTING IT TOGETHER

Thus, to implement an initiating function that supports the universal model, the following modifications are required:

```

template <class Buffers, class Handler>
typename async_result<
    typename handler_type<Handler,
        void(error_code, size_t)>::type>::type } Deduce the initiating function's return type
socket::async_receive(Buffers b, Handler&& handler)
{
    typename handler_type<Handler,
        void(error_code, size_t)>::type
        real_handler(std::forward<Handler>(handler)); } Construct real handler object

    async_result<decltype(real_handler)> } Link initiating function's return value to handler
        result(real_handler);

    ...

    return result.get(); — Return the initiation function's result
}

```

To illustrate how this operates in practice, let us manually work through the steps that the compiler performs for us when we write:

```
auto fut = socket.async_receive(buffers, use_future);
```

First, after expanding uses of the `handler_type` trait to be the “real” handler type, we get:

```

template <class Buffers>
typename async_result<promise_handler<size_t>>::type
socket::async_receive(Buffers b, use_future_t&& handler)
{
    promise_handler<size_t>
        real_handler(std::forward<Handler>(handler));

    async_result<decltype(real_handler)>
        result(real_handler);

    ...

    return result.get();
}

```

Second, we expand the uses of the `async_result` template to get:

```
template <class Buffers>
std::future<size_t>
socket::async_receive(Buffers b, use_future_t&& handler)
{
    promise_handler<size_t>
        real_handler(std::forward<Handler>(handler));

    future<size_t> f = real_handler.p_.get_future();

    ...

    return std::move(f);
}
```

#### 8.4 WHAT HAPPENS WITH PLAIN OL' CALLBACKS

The default implementations of the `handler_type` and `async_result` templates is as follows:

```
template <typename Handler, typename Signature>
struct handler_type {
    typedef Handler type;
};

template <typename Handler>
class async_result {
public:
    typedef void type;
    explicit async_result(Handler& h) { /* No-op */ }
    type get() { /* No-op */ }
};
```

These defaults are used when passing a simple callback (i.e. function object) as a handler. In this case the compiler expands the templates such that the code is effectively:

```
template <class Buffers, typename Handler>
void
socket::async_receive(Buffers b, Handler&& handler)
{
    Handler real_handler(std::forward<Handler>(handler));

    /* No-op */

    ...

    /* No-op */
}
```

This is equivalent to the code in Boost versions 1.53 and earlier. This means that the universal model does not introduce any additional runtime overhead for programs that use callbacks.



## 9 IMPACT ON THE STANDARD

Support for a universal model of asynchronous operations is based entirely on library extensions and does not require any language features beyond those that are already available in C++11.

Adding support for a universal model would involve the addition of two new type traits: `handler_type` and `async_result`. Furthermore, the standard may provide guidance on the use of these type traits in the implementation of asynchronous operations.

The standard library may also be extended to provide seamless support for `std::future` under the universal model. This support is orthogonal to other proposed modifications to `std::future`, such as the addition of the `then()` member function.

If the universal model is adopted then other facilities, such as library-based coroutines, may be considered as separate proposals.

## 10 CONCLUSION

Asynchronous operations have gained widespread use and acceptance in domains such as network programming. In many of these use cases, performance is important and the inherent runtime penalties of `std::future` make it an inappropriate choice for a fundamental building block.

With the universal asynchronous model presented in this paper, we have the ability to select an asynchronous approach that is appropriate to each use case. The universal model supports:

- Callbacks, where minimal runtime penalty is desirable.
- Futures, and not just `std::future` but also future classes supplied by other libraries.
- Coroutines or resumable functions, without adding new keywords to the language.

Perhaps most importantly, with the customisation points that the universal model provides, it can be extended to support other tools for managing asynchronous operations, including ones that we have not thought of yet. The universal model broadens the vocabulary that is available for managing asynchronous control flow, but with a relatively small impact on the standard.

## 11 ACKNOWLEDGEMENTS

The author would like to thank Jamie Allsop for providing extensive feedback, suggestions and corrections. The author also gratefully acknowledges the assistance of Oliver Kowalke in allowing universal model support to be added to his proposed Boost.Fiber library.