

High-Level Effect Handlers in C++

DAN GHICA, Huawei Central Software Institute, Edinburgh, UK

SAM LINDLEY, The University of Edinburgh, UK

MARCOS MAROÑAS BRAVO, Huawei Central Software Institute, Edinburgh, UK

MACIEJ PIROG, Huawei Central Software Institute, Edinburgh, UK

Effect handlers allow the programmer to implement computational effects, such as custom error handling, various forms of lightweight concurrency, and dynamic binding, inside the programming language. We introduce `cpp-effects`, a C++ library for effect handlers with a typed high-level, object-oriented interface. We demonstrate that effect handlers can be successfully applied in imperative systems programming languages with manual memory management. Through a collection of examples, we explore how to program effectively with effect handlers in C++, discuss the intricacies and challenges of the implementation, and show that despite its limitations, `cpp-effects` performance is competitive and in some cases even outperforms state-of-the-art approaches such as C++20 coroutines and the `libmprompt` library for multiprompt delimited control.

1 INTRODUCTION

Effect handlers [Plotkin and Pretnar 2009, 2013] are an expressive control mechanism that allows programmers to define and manage bespoke, fit-for-purpose computational effects. Typical examples include customised error handling, mutable state, input/output, lightweight concurrency, dependency injection, and dynamic binding. Handlers also allow for transparent composition of effects — both with each other and with native effects built into the language. From an engineering methodology perspective, the advantage of handlers is the explicit separation of effect definitions from their programming interface, a collection of *commands* (also known elsewhere in the literature as *operations*). This makes the abstraction ergonomic, as it does not require any exotic conventions to use an effect, in contrast to programming with monads, for instance. It also makes the job of instrumenting existing code possible without extensive rewriting.

One application of effect handlers which is emerging as increasingly important is lightweight cooperative concurrency, that is, concurrency in which all tasks are realised on a single OS thread, without pre-emption. Effect handlers can lead to streamlined and efficient implementations of lightweight (green) threads with different kinds of schedulers, fibers, generators, `async/await`, message-passing actors, etc. Multiple such user-defined concurrency abstractions compose seamlessly, to obtain, for example, nested generators, or threads in which every thread has its own scheduled pool of message-passing actors. The advantage of using handlers is that the programmer can fine-tune all the details (schedulers, communication channels, cancellation policy, error handling), while being able to define them at a relatively high level in a type-safe manner, which would be hard to achieve with lower-level tools, such as a bare context-switching mechanism.

The origins of effect handlers [Plotkin and Pretnar 2009, 2013] lie in the realm of functional programming [Kammar et al. 2013], first as a tool to extend the capabilities of the formal semantics of algebraic effects [Plotkin and Power 2001, 2002, 2003], and later as a distinctive feature of several experimental programming languages [Bauer and Pretnar 2015; Biernacki et al. 2019; Convent et al. 2020; Hillerström and Lindley 2016; Leijen 2017b]. Now more mainstream programming languages are beginning to adopt them as built-in features, for example, Multicore OCaml [Sivaramakrishnan et al. 2021] (soon to be released as part of OCaml 5.0) and Uber’s Pyro language for probabilistic programming [Bingham et al. 2019]. Effect handlers heavily influenced the design of the React GUI

framework [Meta 2022], for instance, inspiring the design of “React Hooks”. Effect handlers are also central to the code navigation functionality of GitHub via the Semantic library [GitHub 2022], which itself depends fundamentally on effect handlers in order to allow analyses to scale modularly to support multiple languages and features.

In this paper, we address the open question of whether effect handlers can be meaningfully exploited in system programming languages, such as C++, which allow for low-level programming with a focus on performance, but also include sophisticated mechanisms for high-level low-cost abstractions, such as classes, templates, and advanced techniques for memory management. We address this problem by implementing and evaluating `cpp-effects`, a C++ library for programming with effect handlers. It is built around the stack-switching mechanism provided by the `boost.context` library [Boost 2022]. Our handlers are one-shot, meaning that suspended computations can only be resumed once. In this style, resumptions are more akin to mutable linear resources, rather than the immutable (and copyable) continuations of functional languages, which is in accord with how stack and memory are manipulated in C++ via the RAII idiom and move semantics [Combette and Munch-Maccagnoni 2018].

The main contributions of this paper are the following:

- (1) **A usable library for effect handlers in C++.** The design goal of our library is to reconcile an ergonomic API with type, memory, and performance guarantees commensurate with existing C++ practice. The library requires no additional compiler support and incurs no performance penalties for code that does not use its features. We argue that this library is not merely a proof-of-concept exercise, but can offer a decisive step towards the incorporation of effect handlers into C++ programming practice.
- (2) **An evaluation of programming with handlers in an object-oriented system programming language.** Effect handlers have been previously implemented in object-oriented languages with automatically managed (garbage-collected) memory [Brachthäuser et al. 2018; Inostroza and van der Storm 2018]. They have also been implemented in C [Leijen 2017a] but without a typed API. We show that abstractions characteristic of system programming (and C++ in particular), such as templates and techniques for direct memory management in the absence of garbage collection, allow for a high-level typed API, and a reasonably efficient implementation. We discuss in depth some further difficulties of programming with effect handlers in such a setting. For example, the style of programming cannot be ported wholesale from functional languages, since such naive implementations will quickly overflow the stack. Such problems require specific solutions, which we provide.
- (3) **Higher-level library on top of a low-level core.** The core `boost.context` library suffers from limited functionality, but enjoys the benefits of small size and wide availability. Our library is built on top of it, in pure C++, so it is as portable as `boost.context` itself. In fact it should be easy to replace `boost.context` with alternative low-level stack manipulation backends. Our design allows for a clear separation between the architecture-specific core and the implementation of handlers using abstractions provided by the language. It can serve as a recipe for libraries in other languages, and in the implementation of compilers of languages with native effect handlers. In the particular case of a compiler, the benefit is that the compiler’s backend need only provide a basic stack-switching mechanism, while the rest of the abstraction can be implemented entirely in the frontend.

The focus of this paper is twofold: the programmer’s experience of using effect handlers in C++ and the implementation of the library.

Section 2 discusses effects and their handlers on some typical examples: mutable state, cooperative threads, message-passing actors, and generators. These examples also illustrate the way effect

handlers can be used in the OO setting: mutable state can be implemented using native mutable state; a scheduler has the familiar shape of a while-loop that picks the process to be resumed; actors can be implemented as a composition of state and threads. We discuss how the lack of tail-call optimisation and the use of value types and templates can adversely affect the experience of using handlers, when compared to Java or functional languages, but we also show how some commonly used workarounds can be successfully applied here. We also discuss how the design of the library allows for some optimisations that can be applied by the programmer to make the code faster and more readable, e.g., a tail- and self-resumptive command clause in a handler can be explicitly marked as such, and in effect can be sped up by avoiding context switching altogether.

Section 3 takes a closer look at the implementation details of effect handlers. Our approach is based on a stack of active handlers, each storing a call-stack segment (a *stacklet*), used to evaluate the handled computation. This is a known implementation strategy, used for example in the Multicore OCaml compiler [Sivaramakrishnan et al. 2021]. However, due to the peculiarities of our particular setting, we must address a number of novel obstacles. One example is reusing the same handler for computations of different types. In existing implementations (functional, OO, and in C), such reuse is realised via reference types, that is, the handler actually manipulates pointers to data. Our library uses templates to allow for value types, that is, direct manipulation of memory, which can be more efficient. This forces us to pay additional attention to the type structure of different elements of the implementation. Another example is managing the lifetime of OO-style handlers, which turns out to be surprisingly subtle. For this reason, the library takes responsibility for allocating and deallocating handlers, but the user is still able to configure this process: for handlers of a certain shape, the user can trade the guaranteed memory safety for performance. The result is a library with explicitly typed commands, handlers, and resumptions, which does not require the user to perform any type casts, but still allows for the convenient dynamic pairing of commands and handlers. (It is still possible to encounter a runtime error when a command does not have a corresponding handler in the context, as we do not provide an effect-tracking system.)

Section 4 presents some microbenchmarks. Though our main focus is the ergonomics of programming with `cpp-effects`, and we are sure there is plenty of scope to improve performance, these benchmarks provide evidence that `cpp-effects` performance can be competitive with the existing mainstream. In particular, we compare with generators built on top of C++20 coroutines, and mutable state built using Section 5 discusses related work and Section 6 concludes.

Our library is available at: <https://github.com/maciejpirog/cpp-effects>.

2 THE PROGRAMMER'S INTERFACE

In this section, we introduce the features of `cpp-effects` through a collection of examples including mutable state, lightweight cooperative threads, actors, and generators. These examples and more (e.g., `async/await` and a toy GUI library inspired by React) are available with the library.

2.1 Mutable state

We begin with an implementation of basic mutable state effect in order to illustrate how effects and handlers work in our library. The main purpose of this example is not to illustrate anything like the full power of effects and handlers, but rather to give a minimal overview of the syntax and semantics of the core features of our library. Having said that, it turns out that mutable state as an effect is a crucial component for our implementation of actors (Section 2.3).

The aim is to abstract over what it means to read and write to a single state cell. The interface will be given by two *commands*: `Put` writes a value to the state cell, and `Get` reads the current value. Different effect handlers provide different implementations of these commands: an implementation

might implement the state cell on the local heap in the standard way, but it might alternatively store it remotely or apply more complex policies such as caching.

Our library occupies the `cpp_effects` namespace, but in this paper we shorten it to `eff`. First let us define the two commands that constitute the effect interface.

```
#include "cpp-effects/cpp-effects.h"
namespace eff = cpp_effects;

template <typename S>
struct Put : eff::command<> {
    S newState;
};

template <typename S>
struct Get : eff::command<S> { };
```

A command is defined as a class (or a struct, which is really just a class whose members are all public) that inherits from the `eff::command` class. The optional template parameter of `command` is the return type for the command. Arguments are supplied to the command as fields of the class. The two commands are defined as templates in order to enable them to be parameterised by the underlying type of the state cell (`s`). The `Put` command takes a single argument `newState` of type `S` and does not return a value. The `Get` command takes no arguments and returns a value of type `S`.

In order to invoke a command we use the `eff::invoke_command` function. For convenience, we define templated wrapper functions for invoking each of our two commands.

```
template <typename S>
void put(S s) {
    eff::invoke_command(Put<S>{{}}, s);
}

template <typename S>
S get() {
    return eff::invoke_command(Get<S>{});
}
```

(Note that the `{}` braces in the initialiser of `Put` is a required trivial initialiser of `Put`'s base class, `eff::command`. Since `command` does not have any member variables, it is always empty, but C++ still requires its initialiser when constructing an object of type `Put`.) We can now write state computations, such as an increment function `inc` for integer state.

```
int inc()
{
    put(get<int>() + 1);
    return get<int>();
}
```

This function increments the current value of the state and returns the updated value of the state.

So far we have given the state interface as a pair of commands and shown how to invoke them. However, we have not yet given them a meaning. In order to do that, we must define a handler. We now define a state handler that simply represents the state cell as a private member.

```
template <typename Answer, typename S>
class Stateful : public eff::handler<Answer, Answer, Put<S>, Get<S>> {
public:
    Stateful(S initialState) : state(initialState) { }
private:
    S state;
```

```

Answer handle_command(Put<S> p, eff::resumption<Answer()> r) override
{
    state = p.newState;
    return std::move(r).tail_resume();
}
Answer handle_command(Get<S>, eff::resumption<Answer(S)> r) override
{
    return std::move(r).tail_resume(state);
}
Answer handle_return(Answer a) override
{
    return a;
}
};

```

A handler is defined as a class that inherits from the `eff::handler` class. The general idea is that a handler handles some computation by interpreting the commands used by that computation and the return value of the computation appropriately. Although not the case here, the return type of the handler and the return type of the computation being handled may differ. These are determined respectively by the return and argument types of the return clause. The template parameters are: the return type of the handler, the return type of the computation being handled, and a list of all handled commands. In this case the final return type and original return type are both `Answer`. We abstract over this answer type and the type of the state cell using a template.

Each command (here `Put` and `Get`) is given a meaning by overloading the `handle_command` virtual member function. The second argument `r` to the `handle_command` function is a *resumption*. It is an object that captures the rest of the computation being handled. Its template parameter reflects the type of the captured computation — it is a function type from the return type of the command to the answer type of the handler. Our use of resumptions in this example is not particularly interesting (we pass the return value to the resumption once at the end of each command clause), but as we shall soon see resumptions are central for allowing effect handlers to express features such as concurrency.

An important aspect of resumptions that is visible here is that they are *one-shot*. They are movable but not copyable. After resuming, the resumption object becomes invalid. When a resumption object goes out of scope the computation it captures is deleted (i.e., the corresponding stack is unwound). Being one-shot allows resumptions to be implemented efficiently using non-copyable data structures such as various kinds of system stack or fibers. Indeed, our implementation is built on the non-copyable fibers of the `boost.context` library.

Both `boost.context` and our library rely on a common C++ idiom known as *move semantics*, which provides a degree of safety to programming with single-use resources, even without a linear type system. The idea is that one can call `resume` only on so-called r-value references, and not expressions that represent actual objects in memory (l-values). This means that one cannot write `res.resume()`, because `res` is an l-value. The sole purpose of `std::move` here is to lift an l-value to an r-value reference, forcing the programmer to include `std::move` in the code, which is supposed to explicitly mark that the programmer gives up the “ownership” of a resource. This approach does not *statically* preclude using the resource twice, e.g. calling `std::move(res).resume()` twice will compile (a form of linear or ownership type system would be able to rule out such cases), but it will lead to a *run time* error.

The `Put` command clause updates the state member with the new value and then invokes the resumption with no argument. The `Get` command clause invokes the resumption with the current

state. The `handle_return` member function defines how to process the final result value. In this case it is simply returned as is.

Having defined a handler we can now invoke it on a computation.

```
int main()
{
    std::cout << eff::handle<Stateful<int, int>>(inc, 100); // Output: 101
}
```

A computation is handled using the function `eff::handle`. The first argument is the computation. Subsequent arguments are forwarded as arguments to the constructor of `Stateful`.

2.2 Cooperative lightweight threads

Now we move onto an example, cooperative lightweight threads, that begins to demonstrate the real power of effect handlers. We begin by defining two commands and convenient wrapper functions.

```
struct Yield : eff::command<> { };
struct Fork : eff::command<> {
    std::function<void()> proc;
};
void yield()
{
    eff::invoke_command(Yield{});
}
void fork(std::function<void()> proc)
{
    eff::invoke_command(Fork{{}, proc});
}
```

The `Yield` command yields control to another lightweight thread. The `Fork` command forks off a new lightweight thread. They can be used in the following code, where `starter` starts a number of threads, each printing out a number in a loop.

```
void worker(int k)
{
    for (int i = 0; i < 10; ++i) {
        std::cout << k;
        yield();
    }
}
void starter()
{
    for (int i = 0; i < 5; ++i) {
        fork(std::bind(worker, i));
    }
}
```

The intention is that calling `starter` in a scheduler will print out a stream of interleaved digits. The `std::bind` function in the body of `starter` yields a zero-argument function (a thunk), which when called applies `worker` to `i` — alternatively, we could use a lambda expression `fork([](){ worker(i); })`.

A handler for `Yield` and `Fork` defines a scheduler. Here we implement a round-robin strategy.

```
using Res = eff::resumption<void()>;
```

```

class Scheduler : public eff::handler<void, void, Yield, Fork> {
public:
    static void Start(std::function<void()> f)
    {
        queue.push_back(eff::wrap<Scheduler>(f));

        while (!queue.empty()) {
            Res resumption = std::move(queue.front());
            queue.pop_front();
            std::move(resumption).resume();
        }
    }
private:
    static std::list<Res> queue;

    void handle_command(Yield, Res r) override
    {
        queue.push_back(std::move(r));
    }

    void handle_command(Fork f, Res r) override
    {
        queue.push_back(std::move(r));
        queue.push_back(eff::wrap<Scheduler>(f.proc));
    }

    void handle_return() override { }
};

```

The scheduler maintains a queue of lightweight threads represented as resumptions. Execution is initiated by the static `start` member function, which is passed a function as an argument, which becomes the body of the first thread. Each thread is created using the function `eff::wrap`, which lifts a function to a resumption by wrapping a handler around its body. For instance, `eff::handle<H>(f)` is equivalent to `eff::wrap<H>(f).resume()`.

The interesting control flow is provided by the command clauses. The `Yield` clause simply pushes the current resumption onto the queue, causing control to return to the scheduling loop. The `Fork` clause again pushes the current resumption onto the queue, but then also places the new forked thread onto the queue too.

In functional programming languages, this example is typically written using a tail-recursive scheduler function [Bauer and Pretnar 2015; Biernacki et al. 2020; Convent et al. 2020; Sivaramakrishnan et al. 2021]. But here we see that it works perfectly well using a loop instead, implementing the familiar pattern of control switching between threads and a central scheduler loop.

One message we would like to convey is that though effect handlers do require a way of capturing resumptions, they do not depend on functional programming and work perfectly well when programming in a primarily imperative style. We can run the scheduler as follows.

```

int main()
{
    Scheduler::Start(starter);
}

```

This will print out 0102103210432104321043210432104321043210432104321432434.

2.3 Actors

One of the key strengths of effect handlers is that they are composable. We now show how to compose our state and lightweight thread handlers in order to implement message-passing actors similar to those of Erlang [Armstrong et al. 1996]. For simplicity, we here give a fixed implementation of actors, but in Section 2.5 we observe that we can treat actors themselves as an effect whose implementation can itself be given by a handler (or indeed different handlers corresponding to different implementation strategies).

Following Erlang, we refer to actors as *processes*. Each process has its own mailbox. Any process can send messages to the mailbox of any other process, but only the owner of a mailbox can read messages from it. We begin by representing a process identifier as a pointer to a queue of messages.

```
using Pid = std::shared_ptr<std::queue<std::any>>;
```

We use the `std::any` type, because mailboxes are heterogeneous, which means that different messages for the same actor can be of different types.

The actor interface is given by the following four functions.

```
Pid spawn(std::function<void()> body);           // spawn a new process and return its process id
Pid self();                                     // return the process id of the current process
template <typename T> void send(Pid p, T msg)    // send msg to p
template <typename T> T receive();              // read a message from my mailbox
```

Here is a simple ping-pong example that uses the interface to spawn a process. The main process then sends a sequence of messages to the child process which just sends the messages back again.

```
void pong()
{
    while (true) {
        auto [pid, n] = receive<std::tuple<Pid, int>>();
        if (n == 0) { return; }
        send<int>(pid, n);
    }
}

void ping()
{
    auto pongPid = spawn(pong);
    for (int i = 1; i <= 10; i++) {
        send<std::tuple<Pid, int>>(pongPid, {self(), i});
        std::cout << receive<int>() << std::endl;
    }
    send<std::tuple<Pid, int>>(pongPid, {self(), 0});
}
}
```

Now we give the implementation of the actor interface, which uses the effects defined in previous sections: state and threads.

```
Pid spawn(std::function<void()> body)
{
    auto mailbox = std::make_shared<std::queue<std::any>>();
    fork( [= ]() {
        eff::handle<Stateful<void, Pid>>(body, mailbox);
    });
    return mailbox;
}
}
```



```

Pid self()
{
    return get<Pid>();
}

template <typename T>
void send(Pid p, T msg)
{
    p->push(msg);
}

template <typename T>
T receive()
{
    auto mailbox = get<Pid>();
    while (mailbox->empty()) { yield(); }
    auto msg = mailbox->front();
    mailbox->pop();
    return std::any_cast<T>(msg);
}

```

The `spawn` function makes use of the `Fork` command to implement the new process as a lightweight thread. It also uses the `Stateful` handler to manage the mailbox in the body of the process. The `self` and `receive` functions access the mailbox using the `Get` command. (It turns out that for this example we only really need a read-only state effect — we don't use `Put` — as the mailbox itself is accessed through a further indirection.) We make use of templates and `std::any_cast` in order to support different message types. If we try to receive a message when the mailbox is empty, then the current process will yield until the mailbox is no longer empty.

We can now run our ping example using the `Scheduler` handler for lightweight threads, with the `Stateful` handler used internally.

```

int main()
{
    Scheduler::Start(std::bind(spawn, ping));
}

```

This outputs the sequence of integers as expected.

A notable aspect of our implementation is that the state effect is bound dynamically as a result of being implemented using handlers. Whenever we call `receive` or `self`, we dynamically bind to the mailbox, so we do not have to carry any information about the current process id around, and the body of the actor looks natural, even if it uses features such as recursion or higher-order functions. Crucially, we could not use a global variable to maintain the state, because each process has its own mailbox.

2.4 Ergonomics

Sometimes the full power of effect handlers is overkill. For instance, in the `Stateful` handler we always invoke the resumption argument in tail position at the end of each command clause. Moreover, the return clause just returns the value it is passed. It is possible to avoid writing some of this boilerplate. Here is a more ergonomic version of `Stateful`.

```

template <typename Answer, typename S>
class Stateful : public eff::flat_handler<Answer, eff::plain<Put<S>>, eff::plain<Get<S>>>> {
public:
    Stateful(S initialState) : state(initialState) { }
}

```

```
private:
    S state;
    void handle_command(Put<S> p) override
    {
        state = p.newState;
    }
    S handle_command(Get<S>) override
    {
        return state;
    }
};
```

The `eff::flat_handler` class automatically inserts an identity return clause — hence it only takes a single `Answer` template argument before the commands. Flat handlers make it easier to define handlers that are truly parametric in the final result type. The problem is due to the C++ template system. For example, in our original definition of the `Stateful` handler, it is not possible to instantiate `Answer` with `void`, as it is used as the argument type of the `handle_return` method.

The *clause modifier* `eff::plain` in the type arguments of `Handler` indicates that a command is a *plain* command, meaning that the resumption given as the argument of `handle_command` must be invoked in tail position at the end of the command clause (we say that it is *tail- and self-resumptive*). As a consequence there is no need to expose the resumption at all to the programmer. The corresponding `handle_command` no longer includes a resumption argument, and the return type is now that of the command instead of the final return type of the handler. As well as improving readability, plain commands also improve performance as there is no need to perform any kind of context switch. So where possible, it is worth marking a command as plain.

Plain clauses are not to be confused with two ways in which we can resume: `resume` (used in the handler for threads) and `tail_resume` (used in the handler for state in Section 2.1). The reason for the latter is the lack of general tail-call optimisation in C++. If we used `return std::move(r).resume()` in the state handler, the `handle_command` frames created when we invoke a command would stay on the call stack until the very end of the handled computation. Since we need to be able to perform any number of commands in a computation, `tail_resume` allows us to avoid stack overflow. It can be used for any resumption as the last handler-related statement in a command clause (or a function called by a command clause), and it can be used to replace the top-most `handle_command` call-frame with the resumed computation. However, from a code-engineering perspective, if the tail-resumed resumption is the one we get as an argument, it makes more sense to use the `eff::plain` modifier, while if the effect juggles a number of resumptions, it is, in our experience, more convenient to trampoline them in a loop as in the lightweight threads example.

The library provides a further clause modifier `eff::no_resume`. This is used to indicate a command clause that never invokes its resumption at all, in other words an exception. As with plain command clauses, the resumption argument is omitted from a `no_resume` command clause. For instance, we might want to extend the lightweight threading interface to support a command to kill the current thread.

```
struct Kill : eff::command<> { };
```

We could then adapt the `Scheduler` handler as follows:

```
class Scheduler : public eff::handler<void, void, Yield, Fork, eff::no_resume<Kill>> {
    // ...
    void handle_command(Kill) override { }
};
```

The main benefit here is clarity. Because of the techniques we use to optimise creation of resumptions (described in Section 3.4), the performance benefit is minimal.

2.5 Actors revisited

In Section 2.3, we gave a fixed implementation of actors. Now we decouple that implementation into an effect interface and a handler, opening up the possibility of easily plugging in alternative implementations. For instance, we might want to swap out the underlying implementation of lightweight threads to use a different scheduler, or we may wish to give a completely different implementation that does not factor through the lightweight threads implementation at all.

The commands and wrapper functions are as follows.

```
using Pid = std::shared_ptr<std::queue<std::any>>;

struct Spawn : eff::command<> {
    std::function<void()> body;
};

struct Self : eff::command<Pid> { };

struct Send : eff::command<> {
    Pid p;
    std::any msg;
};

struct Receive : eff::command<std::any> { };

Pid spawn(std::function<void()> body)
{
    return eff::invoke_command(Spawn({}, body));
}

Pid self()
{
    return eff::invoke_command(Self{});
}

template <typename T>
void send(Pid p, T msg)
{
    eff::invoke_command(Send({}, p, msg));
}

template <typename T>
T receive()
{
    return std::any_cast<T>(eff::invoke_command(Receive{}));
}
```

The commands themselves make use of `std::any`, but the wrapper functions for sending and receiving are template functions and, in particular, the `receive` wrapper performs the type cast.

Now we can move the actual implementation of the commands into a handler.

```
template <typename Answer>
class Act : public eff::flat_handler<Answer, eff::plain<Spawn>, eff::plain<Self>,
    eff::plain<Send>, eff::plain<Receive>> {
    Pid handle_command(Self) override
    {
        return get<Pid>();
    }
};
```

```

}
Pid handle_command(Spawn s) override
{
    auto mailbox = std::make_shared<std::queue<std::any>>();
    fork( [= ]() {
        eff::handle<Stateful<void, Pid>>(s.body, mailbox);
    });
    return mailbox;
}
void handle_command(Send s) override
{
    s.p->push(s.msg);
}
std::any handle_command(Receive) override
{
    auto mailbox = get<Pid>();
    while (mailbox->empty()) { yield(); }
    auto msg = mailbox->front();
    mailbox->pop();
    return msg;
}
}

```

All of the commands are plain and the handler is a parametric flat handler. Now if we want to change the implementation, we can just define a different handler to use in place of `Act`.

This example illustrates a more general problem that we encounter with defining an API for effect handlers as a library: commands cannot be handled polymorphically. This is because if we wanted to handle a command such as

```
template <typename T> struct Receive<T> { };
```

polymorphically in τ , then the corresponding command clauses for the handler would have to be both virtual and templates, which is illegal in C++.

The solution, which is folklore and applicable in many different scenarios, is to split the command into the monomorphic “core” and a polymorphic wrapper. This, however, often requires other types to be factorised in this way. For example, here is a snippet from our implementation of `async/await`:

```

struct GenericFuture {
    std::vector<eff::resumption<void()>> awaiting;
};

template <typename T>
class Future : public GenericFuture {
    std::optional<T> value;
    ...
};

struct Await : eff::command<> {
    GenericFuture* future;
};

template <typename T>
T await(Future<T>* future)
{
    if (*future) { return *(future->value); }
    eff::invoke_command(Await{{}, future}); // Suspend until the value is ready
}

```

```

    return future->Value();
}
template <typename T>
class Scheduler : public eff::handler<void, T, Yield, Await> {
    void handle_command(Await f, Res r) override
    {
        f.future->awaiting.push_back(std::move(r));
        // going back to the scheduler loop...
    }
    // ...
};

```

In this example, we allow the user to create and await a future of any type, but we cannot make `Await` a template. Instead, we decouple the behaviour from the data. Though the `Await` command itself is monomorphic, including a pointer to the `GenericFuture` base class as a member, we wrap its invocation using a template function. The handler itself does not need to know what the concrete type of the future is, just that it inherits from `GenericFuture`.

2.6 Generators

Generators provide a convenient interface for producing a stream of results. They can be implemented using effect handlers in such a way that the user only sees the generator interface and is not exposed to any underlying commands or handlers. Because the type of the handler used internally is statically known, this provides an opportunity for performance gains.

The interface is given by a single `Yield` command which yields a value to the caller (not to be confused with the `Yield` command used by our earlier lightweight threads examples).

```

template <typename T>
struct Yield : eff::command<> {
    T value;
};

```

(Note that there is no issue in defining `Yield` with a template parameter, as long as any handler for `Yield` fixes `T`, and hence handles it monomorphically.) Our wrapper function takes an additional `label` argument. This argument is passed to an overloaded variant of `eff::invoke_command`. It identifies the handler that will be used to handle the command.

```

template <typename T>
void yield(int64_t label, T x)
{
    eff::static_invoke_command(label, Yield<T>{{}, x});
}

```

By default, and for all the examples we have seen up to now, the handler is chosen as the inner-most one that supports the invoked command. This requires dynamically checking runtime type information (RTTI) stored in the handler's `vtable`: note that when invoking a command, we statically know its type, but we have to dynamically inspect the types of active handlers. However, RTTI has a performance penalty, which can often be avoided. By using `eff::static_invoke_command` we are asserting that we know the exact type of the handler associated with `label` and in particular that it handles the `Yield` command with the correct type. This is a potentially dangerous assumption, but offers significant performance benefits, and can be done relatively safely if, as in this case, we encapsulate the command and handler inside a library, so that the user of the library will never need to interact with either.

We now give an implementation of a handler for generators.

```

template <typename T>
struct GenState;

template <typename T>
using Result = std::optional<GenState<T>>;

template <typename T>
struct GenState {
    T value;
    eff::resumption<Result<T>()> resumption;
};

template <typename T>
class GeneratorHandler : public eff::handler<Result<T>, void, Yield<T>> {
    Result<T> handle_command(Yield<T> y, eff::resumption<Result<T>()> r) override
    {
        return GenState<T>{y.value, std::move(r)};
    }
    Result<T> handle_return() override
    {
        return {};
    }
};

```

A generator state comprises a value and a current resumption. The return value of the handler is a `Result<T>` object, which is an optional generator state.

The more interesting part of the implementation is in a separate `Generator` class.

```

template <typename T>
class Generator {
public:
    Generator(std::function<void(std::function<void(T)>> f)>
    {
        auto label = eff::fresh_label();
        result = eff::handle<GeneratorHandler<T>>(label, [f, label](){
            f([label](T x) { yield<T>(label, x); });
        });
    }

    Generator() { } // A dummy generator that generates nothing

    T Value() const
    {
        if (!result) { throw std::out_of_range("Generator::Value"); }
        return result.value().value;
    }

    bool Next()
    {
        if (!result) { throw std::out_of_range("Generator::Next"); }
        result = std::move(result->resumption).resume();
        return result.has_value();
    }

    explicit operator bool() const
    {

```

```

    return result.has_value();
}

```

private:

```

    Result<T> result = {};
};

```

The idea is that when a generator is created its body is executed until the first value is yielded. The constructor takes the body as a higher-order function parameterised by a yield function whose implementation is supplied by the constructor itself. This implementation invokes the `yield` command using a fresh label that identifies the generator handler. (The library maintains a global counter of labels and the user can create a fresh label using `eff::fresh_label1`.) This label is also passed to an overloaded version of `eff::handle` in order to associate it with the handler. Notice that generators are not copyable, since `GenState` is not copyable, since `eff::resumption` is not copyable.

The `Value` member function returns the current value and the `Next` member function moves onto the next value by invoking the resumption, returning `true` if the stream of values has not been exhausted. The following example illustrates how to use a generator to output a stream of 100 natural numbers.

```

int main()
{
    Generator<int> naturals([](auto yield) {
        int i = 1;
        while (true) { yield(i++); }
    });

    for (int i = 0; i < 100; i++) {
        std::cout << naturals.Value() << std::endl;
        naturals.Next();
    }
}

```

Again notice that this user code makes no reference to any commands or effect handlers.

3 IMPLEMENTATION

In this section, we give an overview of the implementation, and detail a few aspects specific to our setting. In general, our approach is based on a stack of handlers, similar to how effect handlers are implemented, for example, in Multicore OCaml [Sivaramakrishnan et al. 2021]. We discuss the inheritance structure of the `handler` class, and memory management of handlers and resumptions. In the code snippets in this section, we frequently assume that the code is inside the namespace `eff`.

3.1 Metastack (the stack of handlers)

We begin by explaining the semantics of effect handlers via manipulation of the call stack, and show how an optimised version of this semantics can directly inform the implementation. Since handlers are a form of generalised resumable exceptions, we first draw a parallel between effect handlers and exception handlers.

The familiar semantics (but not necessarily implementation) of exceptions can be described as follows. The `try comp catch(const E& f) catchClause` statement pushes a handler frame (which corresponds to `catch(const E& f) catchClause`) onto the stack, and proceeds with `comp`. Execution of the `throw e` statement, assuming `e` is of type `E`, unwinds the stack until it finds a frame of the shape that corresponds to some `catch(const E& f) catchClause`, and then continues with `catchClause` with the value `e` bound to the reference `f`. Figure 1 depicts the process of throwing an exception.

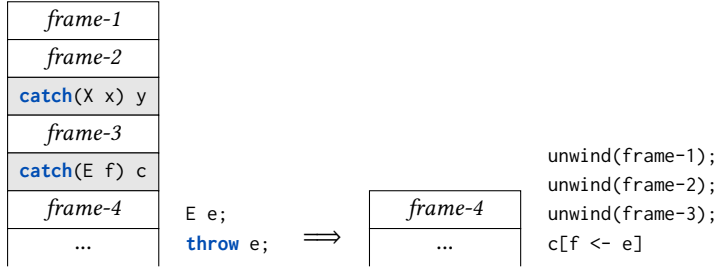


Fig. 1. Unwinding the call-stack on throwing an exception

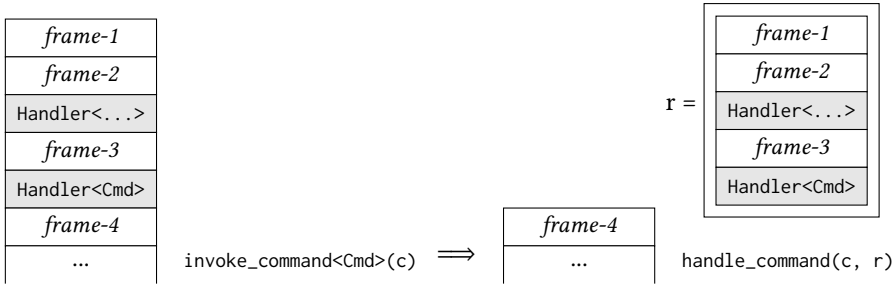


Fig. 2. Invoking a command

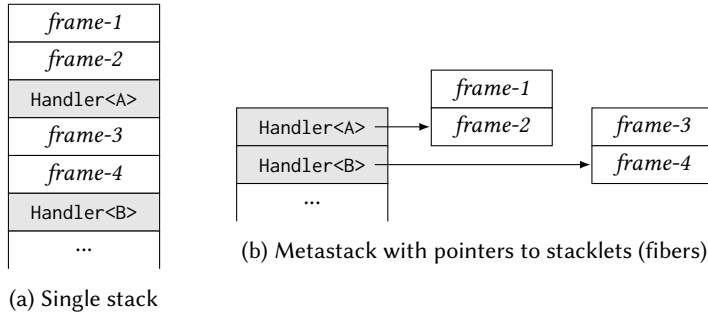


Fig. 3. Different representations of stack

In the case of effect handlers, the function $\text{eff}::\text{handle}\langle H \rangle(f)$ pushes a new frame, which corresponds to a new object of type H , and proceeds with $f()$. A call to $\text{eff}::\text{invoke_command}\langle \text{Cmd} \rangle(c)$ finds the first frame on the stack that corresponds to a handler that can interpret Cmd , and then continues with its $\text{handle_command}\langle \text{Cmd}, \text{eff}::\text{resumption}\langle \dots \rangle \rangle$, which receives as arguments the value c and a new resumption that stores the stack segment above and including the handler frame. Importantly, we do not unwind the stack. The stack stored in a resumption is unwound only when the resumption is deleted (that is, goes out of scope). Figure 2 depicts the invocation of a command.

The stack used in this semantics of effect handlers consists of segments of regular call frames, separated by handler frames. In practice, we can keep each segment of the stack (a *stacklet*) in a separate chunk of memory, and use a *metastack*: a stack of handlers, each with a pointer to the

segment of regular frames above. This representation has two advantages over a flat representation. First, it is faster to search the metastack for the right handler, as we do not have to traverse every frame on the stack. Second, we need not physically move the memory that contains regular frames when we create or resume a resumption; we need only manipulate pointers to stacklets.

3.2 Commands and handlers

The informal semantics given above is the basis of the implementation of the library. It does not require any additional support from the compiler, since the metastack is simply a global data structure (a linked list of pointers to handler objects), while the `boost.context` library provides the mechanism for stacklet allocation and switching.

There are three main class templates in the library: `command`, `handler`, and `resumption`. The `command` template, used as a parent class for user-defined commands, does not provide any functionality of its own. It just specifies the return type of the command and the corresponding resumption type. There are two specialisations of `command`, depending on the number of supplied type parameters:

```
template <typename... Outs>
struct command;

template <typename Out> // Specialisation for a single type parameter
struct command<Out> {
    using out_type = Out;
    template <typename Answer> using resumption_type = Answer(Out);
};

template <> // Specialisation for no type parameters
struct command<> {
    using out_type = void;
    template <typename Answer> using resumption_type = Answer();
};
```

In both cases `out_type` is the return type of the command (the return type of `eff::invoke_command<Cmd>`) and `resumption_type` is the type of resumptions created when a computation is suspended on this command. A command can be handled by different handlers with different `Answer` types, hence `resumption_type` is a template instantiated in the `handle_command` member function of a handler. We require two separate specialisations of `command`, instead of simply defaulting `Out` to `void`, because `Answer(void)` is not a well-formed type in C++.

Handlers are defined via multiple inheritance to implement functionality that allows them to serve as frames on the metastack and provide interpretations to commands. Every handler inherits from the `metaframe` class, which groups a pointer to the stacklet (a *fiber* in `boost.context`'s terminology) and a label that can be used to select a handler as in the generator example in Section 2.6.

```
class metaframe {
    virtual ~metaframe() { }
    boost::context::fiber fiber;
    int64_t label;
    // ...
};
```

The virtual destructor allows us to use runtime type information to implement the (optional) dynamic pairing of commands and handlers described in Section 3.3.

For each command listed in the template arguments pack, `handler` inherits from the `command_clause` class. It provides the virtual member function `handle_command` that particular implementations of handlers must override.

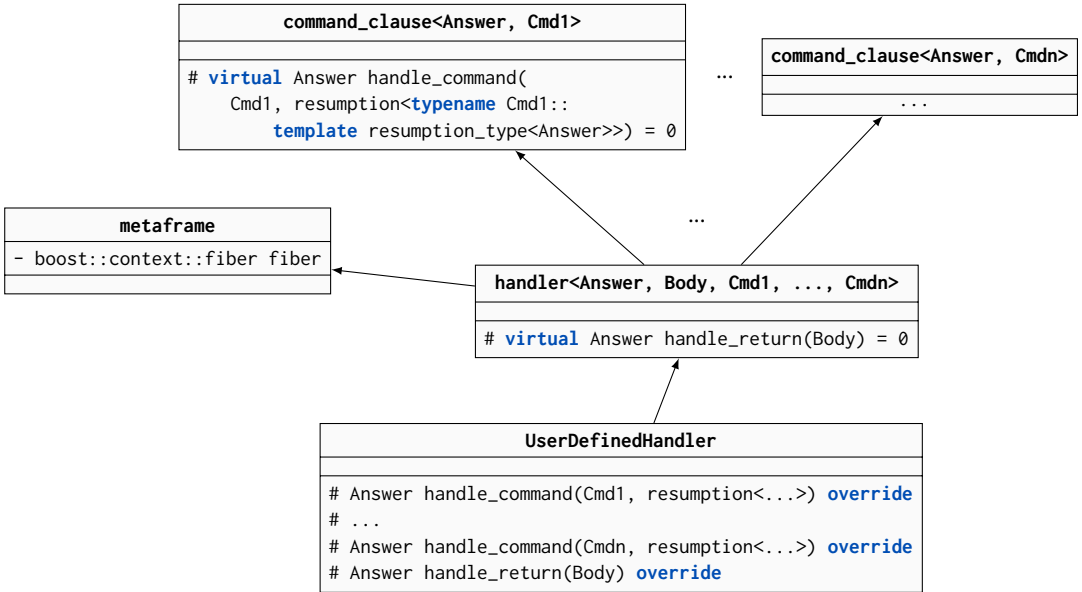


Fig. 4. API for defining handlers

```

template <typename Answer, typename Cmd>
class command_clause {
protected:
    virtual Answer handle_command(Cmd, resumption<typename Cmd::template resumption_type<Answer>>)
        = 0;
    // ...
};

```

Then, the handler is defined as:

```

template <typename Answer, typename Body, typename... Cmds>
class handler : public metaframe, public command_clause<Answer, Cmds>... {
    using command_clause<Answer, Cmds>::handle_command...;
protected:
    virtual Answer handle_return(Body b) = 0;
    // ...
};

```

The `using` declaration in the definition of `handler` exposes `handle_command` from every `command_clause` base, in effect combining them together into one overloaded function. The API of handlers is summarised in Figure 4.

The metastack is a list of pointers to objects of the common superclass of all handlers, `metaframe`. This allows us, for example, to access the label of a handler on the metastack without knowing its actual type. Note that `handler` is a template, so it would be cumbersome, if not impossible, to make the metastack a list of (well-typed pointers to) `handler` objects.

```
std::list<std::shared_ptr<metaframe>> metastack;
```

The fact that the metastack is implemented as a linked list means that we can easily move parts of the metastack to resumptions when invoking a command (as in Figure 2) and back when resuming. Alternatively, we could use a vector, which provides faster handler lookup, but requires allocation

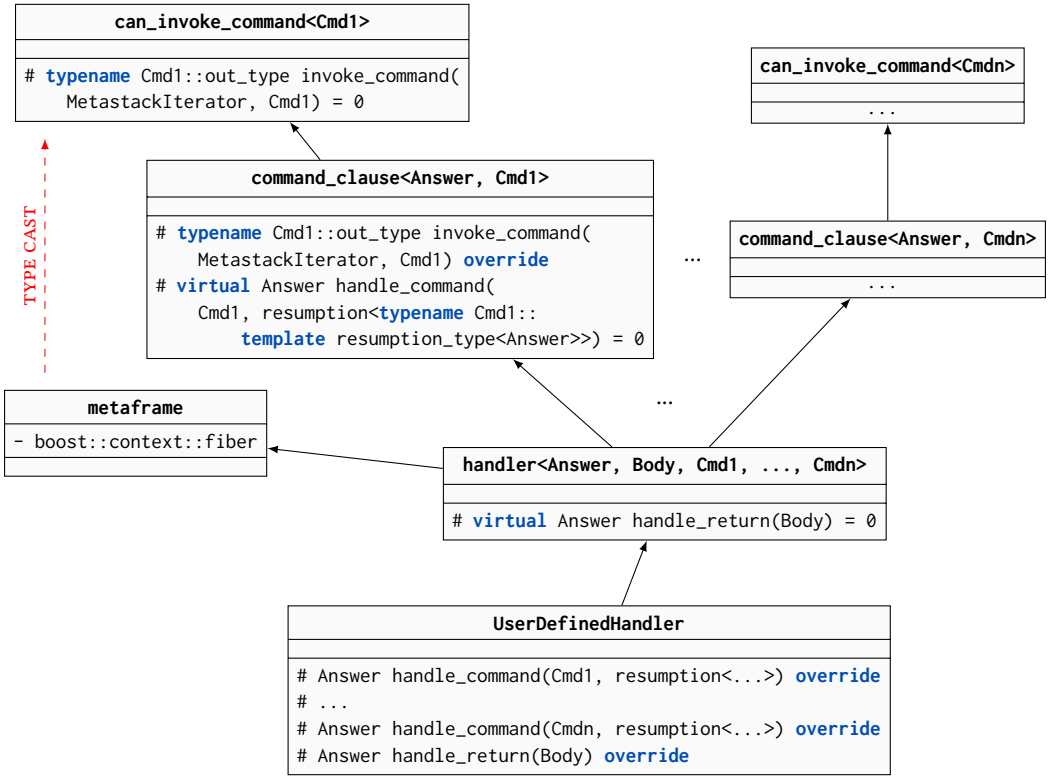


Fig. 5. Internals of handlers needed to invoke commands

every time we create a resumption. In practice, a metastack usually contains no more than a handful of frames, and our experiments suggest that performance of the two possible implementations is similar. We detail why we need reference-counting shared pointers in Section 3.5.

The expression `eff::handle<H>(comp)` first creates a new object of class `H`, pushes a pointer to it on the metastack, and uses its `metaframe::fiber` to evaluate `comp()`. When `comp()` returns, the result is given to the `handle_return` of the handler, which produces the final answer of the call to `eff::handle`.

3.3 Invoking and handling a command

When the user invokes a command `Cmd`, we first need to find the appropriate matching handler on the metastack, as shown in Figure 2. When looking for a handler dynamically, using RTTI, we encounter a problem, which stems from the fact that we have to deal with value types in C++: we cannot dynamically check if a value on the stack inherits from `command_clause<Answer, Cmd>`, simply because in general we cannot know the `Answer` template parameter.

On the other hand, it is easy to find a handler by its label, as we need only compare with the `label` member of the frames on the stack. But then we encounter the same problem when we want to call the `handle_command` function with appropriate arguments later on: we would have to cast the found `metaframe` to (a class that inherits from) `command_clause<Answer, Cmd>`, while we cannot know in general what the `Answer` type is.

We solve this problem by observing that `Answer` is not needed at the site of the command invocation. Thus, we implement the entire logic (creating the resumption, switching stacks, and calling the

virtual function `handle_command`) inside the class `command_clause` where `Answer` is available, but trigger it by making `command_clause` implement the following interface that does not depend on `Answer`.

```
template <typename Cmd>
class can_invoke_command {
protected:
    virtual typename Cmd::out_type invoke_command(
        std::list<std::shared_ptr<metaframe>>::iterator it, const Cmd& cmd) = 0;
};

template <typename Answer, typename Cmd>
class command_clause : public can_invoke_command<Cmd> {
protected:
    virtual typename Cmd::out_type invoke_command(
        std::list<std::shared_ptr<metaframe>>::iterator it, const Cmd& cmd) final override;

    virtual Answer handle_command(Cmd, resumption<typename Cmd::template resumption_type<Answer>>)
        = 0;

    // ...
};
```

The `can_invoke_command<Cmd>::invoke_command` function takes as its arguments the command and an iterator pointing below the found `metaframe`. The handler for `Cmd` can be found on the metastack using RTTI by trying to dynamically cast to `can_invoke_command<Cmd>`. The complete diagram of inheritance for a user-defined handler is shown in Figure 5.

Thus, the overall process of invoking a command is as follows.

- (1) The user calls `eff::invoke_command<Cmd>`, which is responsible for finding the right handler on the metastack. Depending on the overload, it does so using runtime type information (the first handler for which the dynamic cast to `can_invoke_command<Cmd>` succeeds) or the handler's label (in which case we use the dynamic cast only on the first handler with the given label). In Figure 5, this cast from `metaframe` to `can_invoke_command<Cmd1>` is indicated by the dashed arrow. The user can also call `eff::static_invoke_command<H, Cmd>`, in which case we find the handler by the label, and statically cast it to `H`.
- (2) The function `eff::invoke_command` calls the found `can_invoke_command<Cmd>::invoke_command` member function, providing it with the command and the `metaframe` below the found handler, since the found handler's command clause is run on the stacklet of the previous handler (see Figure 2). We create the resumption by moving the top segment of the metastack, switch the stack, and run the result's `handle_command`. If the stack is ever switched back to this place, we return the value with which the resumption was resumed.

The virtual function call to `can_invoke_command<Cmd>::invoke_command` (implemented in the derived class `command_clause<Answer, Cmd>`) solves the problem of the unknown `Answer` type of the handler, but it is also useful for implementing the optimisations provided by clause modifiers. Each clause modifier is implemented as a specialisation of `command_clause`. For example, the implementation of `invoke_command` in `command_clause<Answer, plain<Cmd>>` simply does not create a resumption, and does not switch stacks, but only calls `handle_command`, temporarily removing the top segment of the metastack in case `handle_command` itself uses effects.

3.4 Resumptions

Resumptions are implemented as smart pointers to the `resumption_data` class, which stores a segment of the metastack and some additional members used for transferring data when resuming (which we omit from this description).

```

template <typename Out, typename Answer>
class resumption_data {
    std::list<std::shared_ptr<metaframe>> stored_metastack;
    // ...
};

template <typename Out, typename Answer>
class resumption {
    resumption_data<Out, Answer>* data;
    // ...
};

```

When the user invokes a command, the computation is suspended, a resumption is created by moving an appropriate segment of the metastack to `resumption_data::stored_metastack`, and the appropriate command clause of a handler is called. In our design goals, we assumed that this process should be as fast as possible. As allocation is relatively expensive, we want to avoid allocating a new resumption `resumption_data` object every time a command is called, and deleting it when it is resumed or goes out of scope.

Fortunately, we observe that since our resumptions are one-shot, the user cannot copy them, and so there can exist at most one resumption per suspended computation at a given time. Hence, we can pre-allocate all resumptions that will ever be needed for a particular handler object: one for each supported command. Hence, all `resumption_data` objects can be kept as data members of the handler class. In particular, each `command_clause<Answer, Cmd>` base provides a storage for a resumption that “hangs” on `Cmd`:

```

template <typename Answer, typename Cmd>
class command_clause : public can_invoke_command<Cmd> {
    // ...
    resumption_data<typename Cmd::out_type, Answer> resumption_buffer;
};

```

A local `resumption` object is created as a pointer to the `resumption_buffer` and given as argument to the command clause, which leads to an interesting circular dependency. A resumption contains a stored metastack (as a list, so via pointers), the first metaframe of which is a pointer to a handler which contains the resumption as a member. This cycle is problematic when a resumption goes out of scope, as its destructor needs to break this cycle first, which means that it is not enough to use `std::unique_ptr` to implement the `resumption` class.

3.5 Lifetime of handlers

Prior implementations of effect handlers in object-oriented settings take advantage of automatic memory management: the object representing a handler is created when calling the equivalent of `eff::handle`, and is deleted automatically by the garbage collector. In C++, there is no default garbage collector, and our library manages its own memory, relieving the user from deleting handlers manually.

When should a handler be deleted? One obvious guess would be: when it is popped from the metastack (that is, after the return clause returns) or when a resumption that holds a pointer to it is deleted (the handlers are one-shot, so there can be only one such resumption). However, consider the following simple example of an effect that logs messages:

```

struct Log : eff::command<> { std::string msg; };

class Logger : public eff::handler<std::string, void, Log> {
public:

```

```

    Logger(std::string separator) : separator(separator) { }
private:
    const std::string separator;
    std::string handle_return() override
    {
        return "";
    }
    std::string handle_command(Log l, eff::resumption<std::string()> r) override
    {
        return l.msg + this->separator + std::move(r).resume();
    }
};

```

Since the order of evaluation of operator arguments is unspecified in C++, the body of the command clause for the command `Log` may well be treated by the compiler as equivalent to:

```

auto&& temp = std::move(r).resume();
return l.msg + this->separator + temp;

```

Now observe that `std::move(r).resume()` resumes the entire computation and returns *after* the return clause returns. This means that the second line (`return l.msg + ...`) happens after the object is deleted, and so the expression `this->separator` tries to access a member of a deleted object! This means that we need to keep the handler alive as long as it is on the metastack (or the metastack is stored in a resumption) or there are live stack frames that can refer to it. Since we have no way to statically determine where and when such frames might live, we manage the handler's lifetime using a shared pointer. It is shared pointers that we actually keep on the metastack, and we create a copy of the pointer for the duration of each command clause.

Bumping up the counter on each entry to a command clause and running the destructor of a shared pointer on exit has a performance penalty even up to about 20% of the time of the invoke-resume cycle (see Section 4.1). It is especially unfortunate, because in most cases, there is no need for this, as most examples meet at least one of the following conditions:

- (a) The handler does not expose the resumption to the outside world. In such a case, we know that all command clauses are run on top of the `eff::handle` frame, and so when the frame is removed from the stack (either the function returns or the frame is unwound), we can safely delete the handler, because we know there are no more frames for command clauses anywhere on the stack or stored metastacks.
- (b) The command clauses do not use the internal state of the handler after `resume`. In this case, we can safely delete the handler even if there are still command clauses running on the current stack or any of the stacklets.

In these cases, the library allows the programmer to avoid paying the performance penalty needed in the general case. It is possible via another clause modifier, `no_manage`, that allows the programmer to trade guaranteed memory-safety for performance. This modifier asserts that at least one of the conditions above is true, and the command clause need not memory-manage the handler. A clause that is marked as `no_manage` will not contribute to the reference count of the handler. If all command clauses in the handler are marked as such, it means that there exist at most two references to the handler object: the pointer on the metastack, and a local variable in `eff::handle`. If (a) happens, the metastack pointer is removed first, but the handler is kept alive by the pointer in `eff::handle`. If (b) happens, it might be the case that the pointer in `eff::handle` is removed first (for example, when a command clause stores the resumption in a global variable, and returns an unrelated value), and the handler is deleted after being popped from the metastack, but this will not

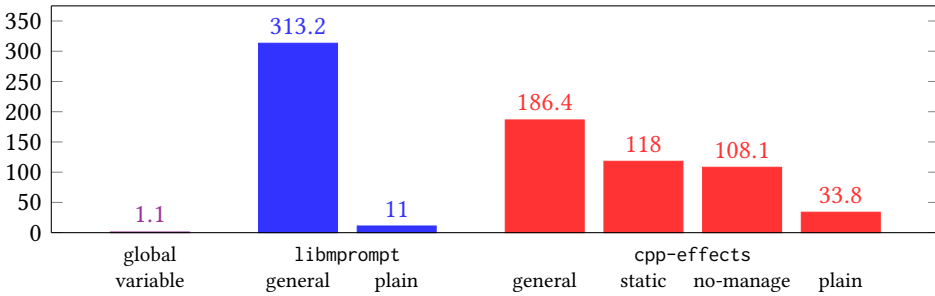


Fig. 6. State benchmark using Clang natively (average time per iteration in ns)

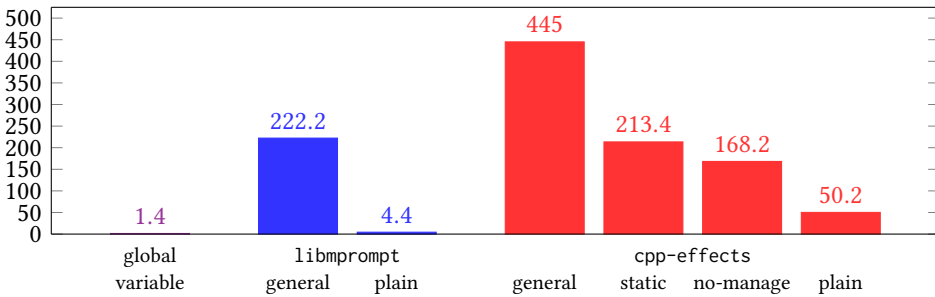


Fig. 7. State benchmark using GCC in Docker (average time per iteration in ns)

cause any problems, as we know that the handled computation has ended (so there will be no new command clauses called) and all live command clause frames have already been called by `resume`.

4 PERFORMANCE

In this section we explore the performance of the library. Since this is, as far as we are aware, the first high-level library for effect handlers in C++, there are no perfect candidates against which to measure performance. Nevertheless, we believe we can still draw some meaningful conclusions about the feasibility of using our library in programming practice.

We compare user-defined effects with those built into the language, such as exceptions and C++20 coroutines. Built-in effects are optimised for specific tasks, so unsurprisingly their performance is markedly better than equivalents implemented via user-defined effects — especially given that we provide *just* a library without any dedicated compiler optimisations. However, for features beyond those provided by the language (for example, resumable exceptions, or stackful coroutines), such a comparison can help to assess the performance penalty of the additional expressiveness.

We also compare `cpp-effects` with the `libmprompt` library [Leijen and Sivamarakrishnan 2022], an existing C implementation of effect handlers. Differences in functionality include that `libmprompt` does not provide a typed high-level interface, features functional-style parameterised handlers instead of OO-style stateful handlers, and supports multi-shot resumptions.

All benchmarks were run on a MacBook Pro Intel Core i5 1.5GHz. We compiled and ran the programs in two set-ups: natively, compiled with Clang 11.0, and in a Docker container running Debian, compiled with GCC 12.1. Interestingly, the relative numbers in the respective cases differ

noticeably. This is because invoking a command (in `cpp-effects` and in `libmprompt`) is a light-weight task, and one can expect that details such as differences in the memory layout provided by the compiler or the peculiarities of the system allocator can significantly affect the running time.

4.1 Mutable state

Our first program consists of a loop that increases the value of a piece of mutable state. The core part of the `cpp-effects` version is thus:

```
eff::handle<Stateful<int64_t>>([]() {
  for (int i = 0; i < 50000000; i++) {
    put(get<int64_t>() + 1);
  }
}, 0);
```

We tested using four different handlers with a gradation of optimisations: a general handler as defined in Section 2.1, a version optimised with static command invocation as in Section 2.6 (recall that “static” refers to knowing statically the type of the handler, which is still searched for dynamically on the stack), a “no-manage” handler (see Section 3.5), and a plain handler (Section 2.4). We also compared to general and “plain” equivalents implemented with `libmprompt`.

The results (average time in nanoseconds of an iteration of the loop) are in Figure 6 for the Clang/native case and in Figure 7 for the GCC/Docker case. One conclusion we can draw is that the dynamic type cast in the general case (Section 3.3) is rather expensive, especially when compiled with GCC. This type cast could be actually avoided if commands were grouped together to form “effects”. This way, each command would determine the rest of the commands in a handler, which in turn would make the vtable layout of the handler always statically known for a given command, and so the dynamic cast could be replaced by a (much cheaper) virtual function call. However, for now we have instead opted for a more readable API. Another conclusion is that the relative performance of `cpp-effects` and `libmprompt` significantly vary with kinds of handlers and the compiler/OS. We suspect that most handlers in practical programming will fall into either the “static, no-manage” category (in which `cpp-effects` seems to be outperforming `libmprompt`), or the “plain” category (in which `libmprompt` definitively outperforms `cpp-effects`).

4.2 C++20 coroutines

We compare generators implemented via effect handlers with coroutines, which were introduced in C++20. The C++ coroutines are *stackless*, which means that they do not run on a separate call-stack, but instead are compiled via a program transformation. In particular, they are compiled to objects with a member function representing the body of the coroutine, and data members representing local variables. Suspending a computation is compiled to returning from the function, and resuming by jumping to a particular place in the function body.

Such jumping in and out of the function body can be much faster than stack switching, which we test on a program that generates consecutive natural numbers in a loop. This generator is resumed a number times, adding up the generated numbers. The results (GCC/Docker only, as Clang 11 does not support coroutines) are shown in Figure 8, relative to an implementation that adds numbers using a loop with no concurrency. Unsurprisingly, native coroutines are much faster than our implementation using effect handlers.

However, the situation is quite different if we benchmark programs that cannot be compiled using a pair of jumps, for example, when the body of the generator is a recursive function. In the case of coroutines, every recursive call corresponds to creating a new coroutine, in effect leading to a heap-allocated stack represented as a linked list of coroutines. Hence, instead of just jumping in

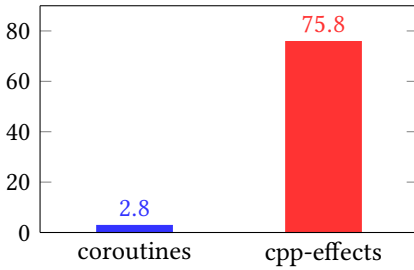


Fig. 8. Generating a number (in ns)

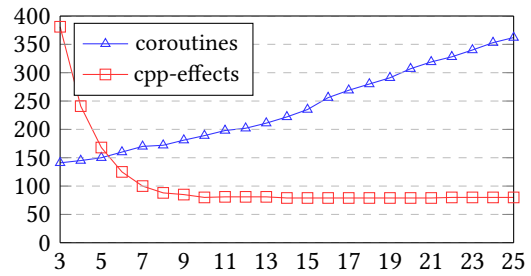


Fig. 9. Recursive tree traversal (ns per node)

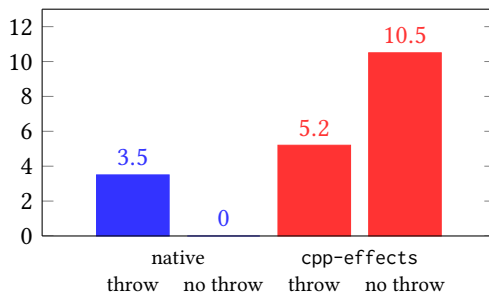


Fig. 10. Exceptions with Clang/native (in ms)

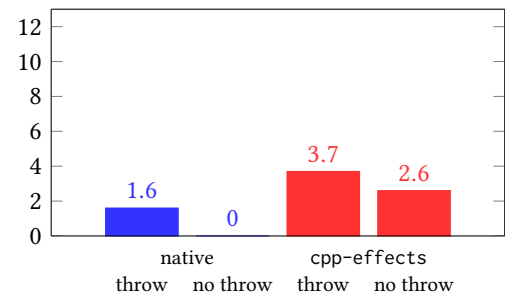


Fig. 11. Exceptions with GCC/Docker (in ms)

and out of a function, the coroutine-based generator must perform many allocations, which has a significant performance cost. We benchmark generators that recursively traverse a full binary tree with values in the leaves. When a generator reaches a leaf, it yields its value. We add up all the values in a tree by repeatedly resuming the generator. The results are shown in Figure 9, where the y -axis is execution time per tree node in nanoseconds, while the x -axis is the depth of the tree. The version based on effect handlers is at least twice as fast for trees of height 8 and more. (Note that the result for each node also amortises the one-off cost of creating the generator, hence the cost per node is relatively high for small trees.)

4.3 Exceptions

Exceptions in idiomatic C++ are primarily used for handling failures which are expected to occur relatively infrequently. This means that most compilers favour the zero-cost implementation, in which installing an exception handler is cheap, while handling an exception is allowed to be expensive, since `throw` is assumed never to be on the hot path of execution. This is quite different to the typical scenario for effect handlers, in which we install an effect handler once (which can therefore be expensive), and then frequently invoke commands (which should therefore be cheap).

Despite the disparity, it is still interesting to compare native exceptions with exceptions implemented using effect handlers. We compare two programs. The first one runs a loop that in each iteration installs a handler and throws an exception. The other one only installs a handler, but does not throw an exception. The results (average time in milliseconds of an iteration of the loop) are shown in Figure 10 for Clang/native and in Figure 11 for GCC/Docker. If no exception is thrown, the native version finishes instantaneously. Interestingly, under Clang/native the version that uses effect handlers and does not throw an exception is slower than the one that throws, whereas the

opposite is the case for GCC/Docker. This reinforces our earlier observation that there can be significant performance disparities between different compiler/OS combinations.

5 RELATED WORK

Effect handlers in C. Leijen [2017a] describes how to implement effect handlers in C on top of `setjmp` and `longjmp`, exposing them via a rather low-level interface using C preprocessor macros. The `libhandler` library [Leijen 2019] is an implementation of this idea. The implementation has to be used with some care, but concrete features such as `async/await` can be implemented in such a way as to expose a relatively safe interface to the programmer. The `libmprompt` library [Leijen and Sivamarakrishnan 2022] is an alternative to the original `libhandler` library. Instead of implementing effect handlers directly it implements `multiprompt` delimited continuations and then a separate library `libmpeff` builds effect handlers on top. The `libmprompt` library uses virtual memory as a way of allowing stacks to grow without ever having to move them.

Effect handlers in OO. JEFF [Inostroza and van der Storm 2018] is an experimental design for an object-oriented programming language with built-in support for effect handlers. The core of JEFF has been formalised and has been implemented in Redex [Klein et al. 2012]. JEFF, unlike `cpp-effects`, supports an effect type system. However, JEFF’s effect type system suffers from similar limitations to Java’s checked exceptions as it does not support effect polymorphism. JEFF, unlike `cpp-effects` which relies on C++-style memory management, assumes a garbage collector. JEFF, in common with `cpp-effects`, takes advantage of dynamic dispatch to define the semantics of effect handlers. However, JEFF hardwires special support for handling effects rather than building all of the functionality on top of existing features as in `cpp-effects`. As such, JEFF also provides mild syntactic sugar for writing effect handlers. Similarly to Koka’s `resume` keyword, which is automatically bound to the resumption of the current handler clause, JEFF binds the current resumption to a special `there` variable somewhat analogous to `this`.

The Java Effekt library [Brachthäuser et al. 2018] is an implementation of effect handlers for Java. Like `libmprompt` it builds effect handlers on top of delimited continuations, which are implemented in Java using a form of CPS translation. The Scala Effekt library [Brachthäuser et al. 2020] is a similar library for Scala, which takes advantage of Scala’s rich type system to incorporate a full-featured effect type system based on capabilities.

Stack switching in WebAssembly. WebAssembly [Rossberg et al. 2018] is a portable low-level bytecode for the web supported by all of the main browser vendors. Work is underway to extend web WebAssembly to support switching between stacks [WebAssembly Community Group 2022] in order to support exactly the kind of features that effect handlers are well-suited for (e.g. `async/await`, lightweight threads, and generators). In particular there is a concrete “Typed Continuations” proposal [Hillerström et al. 2022] along with an implementation in the WebAssembly reference interpreter, which amounts to an extension of WebAssembly with effect handlers. Up to now WebAssembly has largely been used for compiling C and C++; as such the Typed Continuations proposal supports an imperative style of effect handling, not dissimilar to `cpp-effects`, in which schedulers may be implemented as loops rather than with recursive functions and tail-calls.

Clause modifiers. Clause modifiers for handling commands `tail-` and `self-resumptively` (`plain`) or as exceptions (`no_resume`), have been used previously, for instance in the Racket library associated with Kammar et al.’s early work on libraries for effect handlers [Kammar et al. 2013] and in Koka.

Commands and handlers as objects. The idea of representing commands and handlers as objects was introduced by Kammar et al. [2013] in their Haskell library for effect handlers. Like us, they use objects to maintain whatever state is necessary in commands and handlers. However, in their

case, unlike ours, these objects are immutable. Similar ideas arise in later work on designs and implementations of effect handlers that make use of capability-passing [Brachthäuser et al. 2020; Zhang and Myers 2019] and evidence-passing [Xie et al. 2020; Xie and Leijen 2021].

6 CONCLUSION

We are not the first to implement effect handlers in an imperative language or the first to implement effect handlers in an object-oriented programming language. However, as far as we know ours is the first implementation of effect handlers specifically for C++. This presents a particular challenge due to the lack of garbage collection in C++. However, we were successfully able to exploit a broad range of C++ features in order to relatively smoothly integrate effect handlers with C++.

The experience of programming with `cpp-effects` is an extension of the regular experience of programming in C++. Commands and handlers are defined as classes, which can be combined with templates to provide a form of parametric polymorphism in the usual manner. Often a library can be implemented using effect handlers in such a way that the user of the library need not know anything at all about effect handlers (as illustrated by our implementation of generators). Though we have not invested huge effort into trying to optimise `cpp-effects`, it seems to offer adequate performance for realistic use-cases, and in some cases it outperforms existing approaches such as the `libmprompt` library and C++20 coroutines.

In future we would like to explore plugging in alternative backends with a view to improving performance. We would also like to explore means for providing some form of effect type system in order to further tame the complexity of programming in the large with effect handlers. Another direction which it would be interesting to explore is support for multishot effect handlers. This would require an alternative implementation mechanism as it would depend on being able to copy resumptions, but it opens up a range of other applications such as backtracking and probabilistic programming.

REFERENCES

- Joe Armstrong, Robert Virding, Claes Wikström, and Mike Williams. 1996. *Concurrent Programming in Erlang, Second Edition*. Prentice Hall International, Hertfordshire, UK.
- Andrej Bauer and Matija Pretnar. 2015. Programming with algebraic effects and handlers. *J. Log. Algebr. Meth. Program.* 84, 1 (2015), 108–123.
- Dariusz Biernacki, Maciej Piróg, Piotr Polesiuk, and Filip Sieczkowski. 2019. Abstracting algebraic effects. *Proc. ACM Program. Lang.* 3, POPL (2019), 6:1–6:28. <https://doi.org/10.1145/3290319>
- Dariusz Biernacki, Maciej Piróg, Piotr Polesiuk, and Filip Sieczkowski. 2020. Binders by day, labels by night: effect instances via lexically scoped handlers. *Proc. ACM Program. Lang.* 4, POPL (2020), 48:1–48:29. <https://doi.org/10.1145/3371116>
- Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul A. Szerlip, Paul Horsfall, and Noah D. Goodman. 2019. Pyro: Deep Universal Probabilistic Programming. *J. Mach. Learn. Res.* 20 (2019), 28:1–6.
- Boost. 2022. Boost.context library. <https://github.com/boostorg/context>.
- Jonathan Immanuel Brachthäuser, Philipp Schuster, and Klaus Ostermann. 2018. Effect handlers for the masses. *Proc. ACM Program. Lang.* 2, OOPSLA (2018), 111:1–111:27.
- Jonathan Immanuel Brachthäuser, Philipp Schuster, and Klaus Ostermann. 2020. Effekt: Capability-passing style for type- and effect-safe, extensible effect handlers in Scala. *J. Funct. Program.* 30 (2020), e8.
- Guillaume Combette and Guillaume Munch-Maccagnoni. 2018. A resource modality for RAIL. In *LOLA 2018: Workshop on Syntax and Semantics of Low-Level Languages*. 1–4.
- Lukas Convent, Sam Lindley, Conor McBride, and Craig McLaughlin. 2020. Doo bee doo bee doo. *J. Funct. Program.* 30 (2020), e9.
- GitHub. 2022. Semantic library. <https://github.com/github/semantic>.
- Daniel Hillerström, Daan Leijen, Sam Lindley, Matija Pretnar, Andreas Rossberg, and KC Sivamarakrishnan. 2022. WebAssembly Typed Continuations Proposal. <https://github.com/effect-handlers/wasm-spec/proposals/continuations/Explainer.md>.
- Daniel Hillerström and Sam Lindley. 2016. Liberating effects with rows and handlers. In *TyDe*.

- Pablo Inostroza and Tijl van der Storm. 2018. JEff: objects for effect. In *Onward! ACM*, 111–124.
- Ohad Kammar, Sam Lindley, and Nicolas Oury. 2013. Handlers in action. In *ICFP*. ACM, 145–158.
- Casey Klein, John Clements, Christos Dimoulas, Carl Eastlund, Matthias Felleisen, Matthew Flatt, Jay A. McCarthy, Jon Raffkind, Sam Tobin-Hochstadt, and Robert Bruce Findler. 2012. Run your research: on the effectiveness of lightweight mechanization. In *POPL*. ACM, 285–296.
- Daan Leijen. 2017a. Implementing Algebraic Effects in C – “Monads for Free in C”. In *APLAS (Lecture Notes in Computer Science, Vol. 10695)*. Springer, 339–363.
- Daan Leijen. 2017b. Type directed compilation of row-typed algebraic effects. In *POPL*. ACM, 486–499.
- Daan Leijen. 2019. libhandler. <https://github.com/koka-lang/libhandler>.
- Daan Leijen and KC Sivamarakrishnan. 2022. libmprompt. <https://github.com/koka-lang/libmprompt>.
- Meta. 2022. React library. <https://reactjs.org/>.
- Gordon D. Plotkin and John Power. 2001. Semantics for Algebraic Operations. *Electr. Notes Theor. Comput. Sci.* 45 (2001), 332–345.
- Gordon D. Plotkin and John Power. 2002. Notions of Computation Determine Monads. In *FoSSaCS (Lecture Notes in Computer Science, Vol. 2303)*. Springer, 342–356.
- Gordon D. Plotkin and John Power. 2003. Algebraic Operations and Generic Effects. *Applied Categorical Structures* 11, 1 (2003), 69–94.
- Gordon D. Plotkin and Matija Pretnar. 2009. Handlers of Algebraic Effects. In *ESOP (Lecture Notes in Computer Science, Vol. 5502)*. Springer, 80–94.
- Gordon D. Plotkin and Matija Pretnar. 2013. Handling Algebraic Effects. *Logical Methods in Computer Science* 9, 4 (2013).
- Andreas Rossberg, Ben L. Titzer, Andreas Haas, Derek L. Schuff, Dan Gohman, Luke Wagner, Alon Zakai, J. F. Bastien, and Michael Holman. 2018. Bringing the web up to speed with WebAssembly. *Commun. ACM* 61, 12 (2018), 107–115.
- KC Sivamarakrishnan, Stephen Dolan, Leo White, Tom Kelly, Sadiq Jaffer, and Anil Madhavapeddy. 2021. Retrofitting effect handlers onto OCaml. In *PLDI*. ACM, 206–221.
- WebAssembly Community Group. 2022. WebAssembly Stack Switching Extension. <https://github.com/WebAssembly/stack-switching>.
- Ningning Xie, Jonathan Immanuel Brachthäuser, Daniel Hillerström, Philipp Schuster, and Daan Leijen. 2020. Effect handlers, evidently. *Proc. ACM Program. Lang.* 4, ICFP (2020), 99:1–99:29.
- Ningning Xie and Daan Leijen. 2021. Generalized evidence passing for effect handlers: efficient compilation of effect handlers to C. *Proc. ACM Program. Lang.* 5, ICFP (2021), 1–30.
- Yizhou Zhang and Andrew C. Myers. 2019. Abstraction-safe effect handlers via tunneling. *Proc. ACM Program. Lang.* 3, POPL (2019), 5:1–5:29.