

Retrofitting Effect Handlers onto OCaml

KC Sivaramakrishnan, Stephen Dolan, Leo White, Tom Kelly, Sadiq Jaffer and Anil Madhavapeddy

Abstract

Effect handlers have been gathering momentum as a mechanism for modular programming with user-defined effects. Effect handlers allow for non-local control flow mechanisms such as generators, `async/await`, lightweight threads and coroutines to be composablely expressed. We present a design and evaluate a full-fledged efficient implementation of effect handlers for OCaml, an industrial-strength multi-paradigm programming language. Our implementation strives to maintain the backwards compatibility and performance profile of existing OCaml code. Retrofitting effect handlers onto OCaml is challenging since OCaml does not currently have any non-local control flow mechanisms other than exceptions. Our implementation of effect handlers for OCaml: (i) imposes negligible overhead on code that does not use effect handlers; (ii) remains compatible with program analysis tools that inspect the stack; and (iii) is efficient for new code that makes use of effect handlers.

1 Introduction

Effect handlers [50] provide a modular foundation for user-defined effects. The key idea is to separate the definition of the effectful operations from their interpretations, which are given by *handlers* of the effects. For example,

```
effect Read : in_channel -> string
```

declares an *effect* `Read`, which is parameterised with an input channel of type `in_channel`, which when *performed* returns a string value. A computation can perform the `Read` effect without knowing how the `Read` effect is implemented. This computation may be enclosed by different handlers that interpret `Read` differently. For example, `Read` may be implemented by performing a blocking read on the input channel or performing the read asynchronously by offloading it to an event loop such as `libuv`. Thanks to the separation of effectful operations from their interpretations, effect handlers enable new approaches to modular programming.

1.1 An effectful calculator

Effect handlers provide a structured form of delimited continuations [16] that enables non-local control-flow mechanisms such as resumable exceptions, coroutines, generators and asynchronous I/O to be composablely expressed. For example, consider the task of handling division-by-zero in a calculator application. Different calculators handle division-by-zero differently. In the iOS 13 Calculator (`ios`), division-by-zero returns `Error`. In the Google search calculator (`gsearch`) `0/0` returns `Error`, but for `n!=0` the `n/0` returns `Infinity` with the same

sign as `n`. How do we capture all these different behaviours without changing the evaluator?

Common Lisp Conditions [56] are an elegant way to solve this problem. Unlike exceptions where the library decides whether an error is fatal or not, conditions allow the client of the library to make that decision. Conditions provide a *restart* mechanism for non-fatal errors to resume execution from the point where the condition was signalled. On division-by-zero, the calculator signals the handler and lets the handler decide the course of action based on the desired semantics. We now implement this using effect handlers in OCaml syntax:

```
1 type sign = P | N
2 type val = Z of int | Error | Inf of sign
3 type exp = Val of val | ... | Div of exp*exp
4 effect DivBy0 : int -> val
5 let rec eval e =
6   match e with
7   | Val v -> v
8   | ...
9   | Div (e1,e2) ->
10    match eval e1, eval e2 with
11    | Z v1, Z v2 ->
12      if v2 = 0 then perform (DivBy0 v1)
13      else Z (v1/v2)
14    | Error, _ -> Error | Inf s, _ -> Inf s
15    | ...
16 let (/) n d = Div (Val (Z n), Val (Z d))
```

Line 4 declares an effect `DivBy0` that is parameterised by an integer, which when performed returns the result of division-by-zero. If the denominator of the division is `0`, then we perform the `DivBy0` effect with the numerator value (line 12) so that the context may decide the return value. We obtain the `ios` semantics by wrapping the evaluator with this handler:

```
let ios e = match eval e with
| v -> v (* value case *)
| effect (DivBy0 _) k -> continue k Error (* effect case *)
```

If the evaluator returns a value, we simply return that. If the evaluator performs the `DivBy0` effect, we catch the effect in the pattern match along with the continuation `k` of the corresponding *perform delimited* by this handler. For `ios`, we *continue* the suspended computation with an `Error` value and the computation `ios(1/0)` returns `Error`. We can get the `gsearch` semantics by enclosing the evaluator in this handler:

```
let gsearch e = match eval e with
| v -> v
| effect (DivBy0 v) k ->
  if v = 0 then continue k Error
  else if v > 0 then continue k (Inf P)
  else continue k (Inf N)
```

Executing `gsearch(1/0)` returns positive infinity (`Inf P`).

1.2 Asynchronous I/O

We can modularly extend the calculator to read user input using the Read effect defined earlier as follows:

```
type exp = ... | Inp
effect Read : in_channel -> string
let rec eval e = match e with
| ...
| Inp -> int_of_string (perform (Read stdin))
```

The Read effect can then be handled by another handler enclosing the division-by-zero handler:

```
let run e = match gsearch e with
| v -> v
| effect (Read chan) k ->
    register_async_input_line (chan, k);
    continue (get_ui_k ()) ()
```

which registers the read operation to be performed asynchronously with the continuation as the callback. Subsequently, it resumes the user-interface continuation so that the application remains responsive. Observe that neither the division-by-zero handler nor the I/O handler needs to be aware of the other.

1.3 Retrofitting effect handlers

There are many research languages and libraries built around effect handlers [4, 8, 13, 29, 39]. Unlike these efforts, our goal is to retrofit effect handlers onto the OCaml programming language, which has been in continuous use for the past 25 years in large codebases including verification tools [5, 14], mission critical software systems [45] and latency sensitive networked applications [44]. OCaml is particularly favoured for its competitive yet predictable performance, with a fast foreign-function interface (FFI). It has excellent compatibility with program analysis tools such as debuggers and profilers that utilise DWARF stack unwind tables [20] to obtain a backtrace. Our primary motivation for adding effect handlers to OCaml is to support: (i) highly concurrent applications such as web servers using a natural *direct-style* as opposed to using callbacks; and (ii) efficient nested parallel programs in the vein of Cilk [9], OpenMP [33], X10 [11], Intel TBB [32], and variants of Parallel ML [25, 55, 59].

OCaml currently does not support any non-local control flow mechanisms other than exceptions. This makes it particularly challenging to implement the delimited continuations necessary for effect handlers without sacrificing the desirable properties of OCaml. One standard way of implementing continuations is to use continuation-passing style (CPS) in the compiler's intermediate representation (IR) [39]. OCaml does not use a CPS IR, and changing the compiler to utilise a CPS IR would be an enormous undertaking that would affect the performance profile of existing OCaml applications due to the increased memory allocations as the continuation closures get allocated on the heap [22]. Moreover, with CPS, stack unwinding is meaningless and hence, we would lose compatibility with tools that inspect the program stack.

Hence, we choose not to use CPS translation and represent the continuations as call stacks.

The search for an expressive effect system that guarantees that all the effects performed in the program are handled (*effect safety*) in the presence of advanced features such as polymorphism, modularity and generativity is an active area of research [6, 7, 29, 39]. We do not focus on this question in this paper, and our implementation of effect handlers in OCaml does not guarantee effect safety. For example, running `perform (Read stdin)` without an enclosing handler raises `Unhandled exception` at the point of `perform`. We leave the question of effect safety for future work.

1.4 Requirements

We motivate our effect handler design based on the following ideal requirements:

R1: Backwards compatibility. A well-behaved OCaml program does not break under OCaml extended with effect handlers. OCaml code that does not use effect handlers will pay negligible performance and memory cost.

R2: Tool compatibility. OCaml programs with effect handlers produce well-formed backtraces and remain compatible with program analysis tools such as debuggers and profilers that inspect the stack using DWARF unwind tables.

R3: Effect handler efficiency. The program must accommodate millions of continuations at the same time to support highly-concurrent applications. Installing effect handlers, capturing and resuming continuations must be fast.

R4: Forwards compatibility. While this work does not present an effect system, the design must be compatible with a future effect system for OCaml. As a cornerstone of modularity, we also want blocking I/O code to *transparently* be made asynchronous with the help of effect handlers.

The need to host millions of continuations at the same time rules out the use of a large contiguous stack space as in C for continuations. Instead, we resort to using small initial stacks and growing the stacks on demand. As a result, OCaml functions, irrespective of whether they use effect handlers, need to perform stack overflow checks, and external C functions (which do not have stack overflow checks) must be performed on a separate system stack. Additionally, we must generate DWARF stack unwind tables for stacks that may be non-contiguous. In this work, we develop the compiler and runtime support for implementing efficient effect handlers for OCaml that satisfy these requirements.

Our work is also timely. The WebAssembly [28] community group is considering effect handlers as one of the mechanisms for supporting concurrency, asynchronous I/O and generators [58]. Project Loom [42] is an OpenJDK project that adds virtual threads and delimited continuations to Java. The Swift roadmap [57] includes *direct-style* asynchronous programming and structured concurrency as milestones. We believe that our design choices will inform similar choices to be made in other industrial-strength languages.

1.5 Contributions

Our contributions are to present:

- the design and implementation of one-shot effect handlers for OCaml. Our design retains OCaml’s compatibility with program analysis tools that inspect the stack using DWARF unwind tables. We have validated our DWARF unwind tables with the assistance of an automated validator tool [3].
- a formal operational semantics for the effect handler implementation in OCaml.
- extensive evaluation which shows that our implementation has negligible impact on code that does not use effect handlers, and serves as an efficient foundation for scalable concurrent programming.

We have implemented effect handlers in a multicore extension of the OCaml programming language which we call **Multicore OCaml** to distinguish it from stock OCaml.

The remainder of the paper continues with a description of the stock OCaml program stack (§2). We then refine the design further by focussing on retrofitting effect handlers into a mainstream *systems* language (§3), followed by the compilation scheme and operational semantics for Multicore OCaml effect handlers (§4). We then discuss the compiler and the runtime system support for implementing effect handlers (§5), and present an extensive performance evaluation of effect handlers (§6) against our design goals (§1.4). Finally, we discuss the related work (§7) and conclude (§8).

2 Background: OCaml Stacks

The main challenge in implementing effect handlers in Multicore OCaml is managing the program stack and preserving its desirable properties. In this section, we provide an overview of the program stack and related mechanisms in stock OCaml.

Consider the layout of the stock OCaml stack for the program shown in Figures 1a and 1b. The OCaml main function `omain` installs two exception handlers `h1` and `h2` to handle the exceptions `E1` and `E2`. `omain` calls the external C function `ocaml_to_c`, which in turn calls back into the OCaml function `c_to_ocaml`, which raises the exception `E1`. OCaml supports raising exceptions in C as well as throwing exceptions across external calls. Hence, the exception `E1` gets caught in the handler `h1`, and `omain` returns 42. The layout of the stack in the native code backend just before raising the exception in `c_to_ocaml` is illustrated in Figure 1c. Note that the stack grows downwards.

OCaml uses the same program stack as C, and hence the stack has alternating sequences of C and OCaml frames. However, unlike C, OCaml does not create pointers into OCaml frames. OCaml uses the hardware support for `call` and `return` instructions for function calls and returns. OCaml does not perform explicit stack overflow checks in code, and, just like C, relies on the guard page at the end of the stack

region to detect stack overflow. Stack overflow is detected by a memory fault and a `Stack_overflow` exception is raised to unwind the stack.

2.1 External calls and callbacks

OCaml does not use the C calling convention. In particular, there are no callee-saved registers in OCaml. In the x86-64 backend, the OCaml runtime makes use of two C callee-saved registers for supporting OCaml execution. The register `r15` holds the *allocation pointer* into the minor heap used for bump pointer allocation, and `r14` holds a reference to the `Cam1_state`, a table of global variables used by the runtime. This makes external calls extremely fast in OCaml. If the external function does not allocate in the OCaml heap, then it can be called directly and no bookkeeping is necessary. For external functions which allocate in the OCaml heap, the cached allocation pointer is saved to `Cam1_state` before the external call and it is restored on return. Similarly, callbacks into OCaml are also cheap: these involve loading the arguments in the right registers and calling the OCaml function. OCaml callbacks are relatively common as the garbage collector (GC) finalisation functions are executed as OCaml callbacks.

2.2 Exception handlers

The lack of callee-saved registers also makes exception handling fast as there are no callee-saved registers to be restored when the exception handler runs. Installing an exception handler simply pushes the program counter (`pc`) of the handler and the current exception pointer (`exn_ptr` – a field in `Cam1_state`) to the stack. After this, the current exception pointer is updated to be the current stack pointer (`rsp`). This creates a linked-list of exception handler frames on the stack as shown in Figure 1c. Raising an exception simply sets `rsp` to `exn_ptr`, loads the saved `exn_ptr`, and jumps to the `pc` of the handler.

In order to forward exceptions across C frames, the C stub function `cam1_call_ocaml`, pushes an exception handler frame that either forwards the exception to the innermost OCaml exception handler (`raise_exn_c` in Figure 1c) or prints a fatal error (`fatal_uncaught`) if there are no enclosing handlers. Exceptions are so cheap in OCaml that it is common to use them for *local* control flow.

2.3 Stack unwinding

OCaml generates *stack maps* in order to accurately identify roots on the stack for assisting the GC. For every call point in the program, the OCaml compiler emits the size of the frame and the set of all live registers in the frame containing pointers into the heap. During a GC, the OCaml stack is walked and the roots are marked, skipping over the C frames.

OCaml also generates precise DWARF unwind information for OCaml, thanks to which debuggers such as `gdb` and

```

331 1 external ocaml_to_c
332 2 : unit -> int = "ocaml_to_c"
333 3 exception E1
334 4 exception E2
335 5 let c_to_ocaml () = raise E1
336 6 let _ = Callback.register
337 7 "c_to_ocaml" c_to_ocaml
338 8 let omain () =
339 9   try (* h1 *)
340 10    try (* h2 *) ocaml_to_c ()
341 11    | with E2 -> 0
342 12    | with E1 -> 42;;
343 13 let _ = assert (omain () = 42)

```

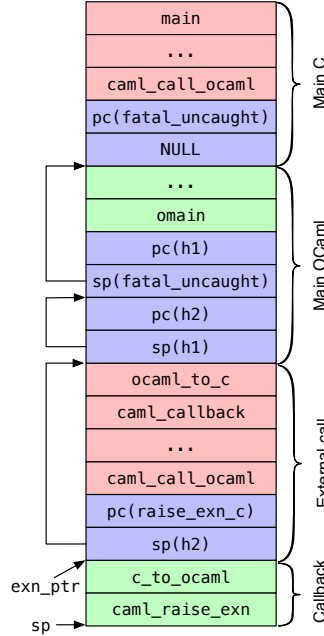
(a) meander.ml

```

344 1 #include <caml/mlvalues.h>
345 2 #include <caml/callback.h>
346 3
347 4 value ocaml_to_c (value unit) {
348 5   caml_callback(*caml_named_value
349 6   ("c_to_ocaml"), Val_unit);
350 7   return Val_int(0);
351 8 }

```

(b) meander.c



(c) Stack layout before raise E1.

```

#0 0x925dc in caml_raise_exn ()
#1 0x6fd3e in camlMeander__c_to_ocaml_83 () at
meander.ml:5
#2 0x925a4 in caml_call_ocaml ()
#3 0x8a84a in caml_callback_exn (...) at
callback.c:145
#4 caml_callback (...) at callback.c:199
#5 0x76e0a in ocaml_to_c (unit=1) at meander.c:5
#6 0x6fd77 in camlMeander__omain_88 () at
meander.ml:10
#7 0x6fe92 in camlMeander__entry () at
meander.ml:13
#8 0x6f719 in caml_program ()
#9 0x925a4 in caml_call_ocaml ()
#10 0x92e4c in caml_startup_common (...) at
startup_nat.c:162
#11 0x92eab in caml_startup_exn (...) at
startup_nat.c:167
#12 caml_startup (...) at startup_nat.c:172
#13 0x6f55c in main (...) at main.c:44

```

(d) gdb backtrace before raise E1.

Figure 1. Program stack on stock OCaml.

lldb, and profilers such as perf work out-of-the-box. For example, for the program in Figures 1a and 1b, one could set a break point in gdb at `caml_raise_exn` to get the backtrace in Figure 1d which corresponds to the stack in Figure 1c.

The same backtrace can also be obtained by using *frame pointers* instead of DWARF unwind tables. OCaml allows compiling code with frame pointers, but they are not enabled by default. The OCaml stack tends to be deep with small frames due to the pervasive use of recursive functions, not all of which are tail-recursive. Hence, the addition of frame pointers can significantly increase the size of the stack¹. Moreover, not using frame pointers saves two instructions in the function prologue and epilogue, and makes an extra register (rbp on x86_64) available. Note that the DWARF unwind information is complementary to the information used by OCaml to walk the stack for GC.

3 Refining the design

In this section, we refine the effect handler design for the OCaml use case. The primary motivation for adding effect handlers to OCaml is to support a non-local control flow mechanism to implement asynchronous I/O in direct style, generators and enable task parallelism. These use cases do not resume the continuations more than once. Hence, our continuations are one-shot, and resuming the continuation more than once raises an `Invalid_argument` exception. It is well-known that one-shot continuations can be implemented efficiently [10].

While OCaml permits throwing exceptions across C frames, we do not allow effects to propagate across C frames as the C frames would become part of the captured continuation. Managing C frames as part of the continuation is a complex endeavour [38], and we find that the complexity budget outweighs the relatively fewer mechanisms enabled by this addition in our setting.

The interaction of non-local control flow with systems programming is quite subtle [19, 40]. Consider the following function that uses blocking I/O functions from the OCaml standard library to copy data from the input channel `ic` to the output channel `oc`:

```

let copy ic oc =
  let rec loop () =
    output_string oc ((input_line ic) ^ "\n");
    loop ()
  in
  try loop () with
  | End_of_file -> close_in ic; close_out oc
  | e -> close_in ic; close_out oc; raise e

```

The function `input_line` raises an `End_of_file` exception on reaching the end of input, which is handled by the exception handler which closes the channels. The `close_*` functions do nothing if the channel is already closed. The code is written in a defensive style to handle other exceptional cases such as the channels being closed externally. Both `input_line` and `output_string` raise `Sys_error` exception if the channel is closed. In this case, the catch-all exception handler closes the channels and reraises the exception to communicate the exceptional behaviour to the caller.

¹<https://github.com/ocaml/ocaml/issues/5721#issuecomment-472965549>

One of our goals (§1.4) is to make this code transparently asynchronous. We can define effects for performing the I/O operations and wrap them up in functions with the same signature as the one from the standard library:

```
effect In_line : in_channel -> string
effect Out_str : out_channel * string -> unit
let input_line ic = perform (In_line ic)
let output_string oc s = perform (Out_str (oc, s))
```

We can then use a top-level I/O handler, similar to the introductory example (§1.2), to discharge the I/O operations asynchronously and resume with the result. This handles value return cases, but what about the exceptional cases `End_of_file` and `Sys_error`? To this end, we introduce a `discontinue` primitive to resume a continuation by raising an exception. In this example, on reaching the end of file, we would discontinue the captured continuation of the `input_line` function with `discontinue k End_of_file`, which raises the exception at `input_line`, and the open channels will be closed.

OCaml programs that use resources such as channels are usually written defensively with the assumption that calling a function will return *exactly once*, either normally or exceptionally. Since effect handlers in Multicore OCaml do not ensure that all the effects are handled, if the function performs an effect with no matching handler, then the function will not return at all. To remedy this, when such an effect bubbles up to the top-level, we discontinue the continuation with an `Unhandled` exception so that the exception handlers may run and clean up the resources.

4 Semantics

In this section, we formalise the effect handler design for Multicore OCaml.

4.1 Static semantics

As mentioned earlier, effect handlers in Multicore OCaml do not guarantee effect safety, but only guarantee type safety. Programs without matching effect handlers are well-typed Multicore OCaml programs. As a result, our static semantics is simpler than languages that ensure effect safety [4, 6, 13, 29, 39, 52]. This is important for backwards compatibility as our goal is to retrofit effect handlers to a language with large legacy codebases; programs that do not use effects remain well-typed, and those that do compose well with those that don't.

The static semantics of effect handlers in OCaml is captured succinctly by its API:

```
type 'a eff = ..
type ('a, 'b) continuation
val perform: 'a eff -> 'a
val continue: ('a, 'b) continuation -> 'a -> 'b
val discontinue: ('a, 'b) continuation -> exn -> 'b
(* Internal API *)
type 'a comp = unit -> 'a
type ('a, 'b) handler =
```

```
{retc: 'a -> 'b;
 effc: 'c. 'c eff -> ('c, 'b) continuation -> 'b; }
val match_with: 'a comp -> ('a, 'b) handler -> 'b
```

We introduce an extensible variant type [49] `'a eff` of effect values, which when performed using the `perform` primitive returns an `'a` value. Constructors for the value of type `'a eff` are declared using the effect declarations. For example, the declaration `effect E : string -> int` is syntactic sugar for adding a new constructor to the variant type `type _ eff += E : string -> int eff`. We introduce the type `('a, 'b) continuation` of delimited continuations which expects a `'a` value for resumption and returns a `'b` value. The continuations may be continued with a suitably typed value or discontinued with an exception.

For handling the effects, our implementation extends OCaml's `match ... with` syntax with effect patterns. The expression

```
match e with
| None -> false | Some b -> b
| effect (E s) k1 -> e1 | effect (F f) k2 -> e2
```

is translated to the equivalent of

```
match_with (fun () -> e)
{ retc = (function None -> false | Some b -> b);
  effc = (function
    | (E s) -> (fun k1 -> e1)
    | (F f) -> (fun k2 -> e2)
    | e -> (fun k -> match perform e with
      | v -> continue k v
      | exception e -> discontinue k e)); }
```

For the sake of exposition, we introduce a `('a, 'b) handler` type. This handler handles a `'a comp` that returns a `'a` value, and itself returns a `'b` value. The handler has a return case `retc` of type `'a -> 'b`. The effect case `effc` handles effects of type `'c eff` with `('c, 'b) continuation` and returns a value of type `'b`. The last case in `effc` *reperforms* any unmatched effect to the outer handler and returns the value and exceptions back to the original performer. In the implementation, *reperform* is implemented as a primitive to avoid executing code on the resumption path.

4.2 Dynamic semantics

We present an operational semantics for a core language of effect handlers that faithfully captures the semantics of the Multicore OCaml implementation. An executable version of the semantics, implemented as an OCaml interpreter, along with examples, is included in the supplementary material.

4.2.1 Syntax. Our expressions (Figure 2a) consist of integer constants (n), variables (x), abstraction ($\lambda x.x$), application ($e e$), arithmetic expressions ($e \odot e$) where \odot ranges over $\{+, -, *, /\}$, raising exceptions (**raise** $l e$), performing effects (**perform** $l e$), and handling effects (**match** e **with** h). Abstractions come in two forms: OCaml abstractions (λ^o) and C abstractions (λ^c). The handler consists of a return case (**return** $x \mapsto e$), zero or more exception cases (**exception** $l x \mapsto$

Constants	$n := \mathbb{Z}$	Handler Closures	$\eta := (h, \epsilon)$
Abstractions	$\Lambda := \lambda^o \mid \lambda^c$	Frame List	$\psi := [] \mid r :: \psi$
Expressions	$e := n \mid x \mid e e \mid \Lambda x. e \mid e \odot e \mid \mathbf{raise} \ l \ e$ $\mid \mathbf{match} \ e \ \mathbf{with} \ h \mid \mathbf{perform} \ l \ e$	Fibers	$\varphi := (\psi, \eta)$
Handlers	$h := \{\mathbf{return} \ x \mapsto e\} \mid \{\mathbf{exception} \ l \ x \mapsto e\} \uplus h$ $\mid \{\mathbf{effect} \ l \ x \ k \mapsto e\} \uplus h$	Continuations	$k := [] \mid \varphi \triangleleft k$
Values	$v := n \mid k \mid \langle \Lambda x. e, \epsilon \rangle \mid \mathbf{eff} \ l \ k \mid \mathbf{exn} \ l$	C stacks	$\gamma := [\psi, \omega]_c$
Frames	$r := \langle e \ \epsilon \rangle_a \mid \langle v \rangle_f \mid \langle \odot e \ \epsilon \rangle_{b1} \mid \langle \odot \mathbb{N} \rangle_{b2}$	OCaml stacks	$\omega := [k, \gamma]_o \mid \bullet$
Environments	$\epsilon := \emptyset \mid \epsilon[x \mapsto v]$	Stacks	$\sigma := \gamma \mid \omega$
		Terms	$\tau := e \mid v$
		Configurations	$\mathfrak{C} := \ \tau, \epsilon, \sigma\ $

(a) Syntax of expressions and configurations

$$\text{STEP C} \quad \frac{(\tau, \epsilon, \psi, \omega) \xrightarrow{c} \mathfrak{C}}{\|\tau, \epsilon, [\psi, \omega]_c\| \rightarrow \mathfrak{C}} \quad \text{STEP O} \quad \frac{(\tau, \epsilon, k, \gamma) \xrightarrow{o} \mathfrak{C}}{\|\tau, \epsilon, [k, \gamma]_o\| \rightarrow \mathfrak{C}}$$

(b) Top-level reductions

VAR	$(x, \epsilon, \psi) \rightsquigarrow (\epsilon(x), \epsilon, \psi)$
ARITH1	$(e_1 \odot e_2, \epsilon, \psi) \rightsquigarrow (e_1, \epsilon, \langle \odot e_2 \ \epsilon \rangle_{b1} :: \psi)$
ARITH2	$(n_1, _, \langle \odot e_2 \ \epsilon \rangle_{b1} :: \psi) \rightsquigarrow (e_2, \epsilon, \langle \odot n_1 \rangle_{b2} :: \psi)$
ARITH3	$(n_2, \epsilon, \langle \odot n_1 \rangle_{b2} :: \psi) \rightsquigarrow (\llbracket n_1 \odot n_2 \rrbracket, \epsilon, \psi)$
APP1	$(e_1 \ e_2, \epsilon, \psi) \rightsquigarrow (e_1, \epsilon, \langle e_2 \ \epsilon \rangle_a :: \psi)$
APP2	$(\Lambda x. e, \epsilon, \psi) \rightsquigarrow (\langle \Lambda x. e, \epsilon \rangle, \epsilon, \psi)$
APP3	$(\langle \Lambda x. e_1, \epsilon_1 \rangle, _, \langle e_2 \ \epsilon_2 \rangle_a :: \psi) \rightsquigarrow (e_2, \epsilon_2, \langle \langle \Lambda x. e_1, \epsilon_1 \rangle \rangle_f :: \psi)$
RESUME1	$(k, _, \langle e_1 \ \epsilon_1 \rangle_a :: \langle e_2 \ \epsilon_2 \rangle_a :: \psi) \rightsquigarrow (e_1, \epsilon_1, \langle k \rangle_f :: \langle e_2 \ \epsilon_2 \rangle_a :: \psi)$
RESUME2	$(\langle \Lambda x. e_1, \epsilon_1 \rangle, _, \langle k \rangle_f :: \langle e_2 \ \epsilon_2 \rangle_a :: \psi) \rightsquigarrow (e_2, \epsilon_2, \langle k \rangle_f :: \langle \langle \Lambda x. e_1, \epsilon_1 \rangle \rangle_f :: \psi)$
PERFORM	$(\mathbf{perform} \ l \ e, \epsilon, \psi) \rightsquigarrow (e, \epsilon, \langle \mathbf{eff} \ l \ [], (\mathbf{return} \ x \mapsto x), \emptyset \rangle_f :: \psi)$
RAISE	$(\mathbf{raise} \ l \ e, \epsilon, \psi) \rightsquigarrow (e, \epsilon, \langle \mathbf{exn} \ l \ \rangle_f :: \psi)$

(c) Administrative Reductions – $(\tau, \epsilon, \psi) \rightsquigarrow (\tau', \epsilon, \psi')$.

ADMINC	$(\tau, \epsilon, \psi, \omega) \xrightarrow{c} \ \tau', \epsilon', [\psi', \omega]_c\ $	if $(\tau, \epsilon, \psi) \rightsquigarrow (\tau', \epsilon', \psi')$
CALLC	$(v, _, \langle \langle \lambda^c x. e, \epsilon \rangle \rangle_f :: \psi, \omega) \xrightarrow{c} \ e, \epsilon[x \mapsto v], [\psi, \omega]_c\ $	
CALLBACK	$(v, _, \langle \langle \lambda^o x. e, \epsilon \rangle \rangle_f :: \psi, \omega) \xrightarrow{c} \ e, \epsilon[x \mapsto v], [k, [\psi, \omega]_c]_o\ $	if $k = [], (\mathbf{return} \ x \mapsto x), \emptyset$
RETTOO	$(v, \epsilon, [], [k, \gamma]_o) \xrightarrow{c} \ v, \epsilon, [k, \gamma]_o\ $	
EXNFWD O	$(v, \epsilon, \langle \mathbf{exn} \ l \ \rangle_f :: _, [\psi, \eta] \triangleleft k, \gamma]_o) \xrightarrow{c} \ v, \epsilon, [(\langle \mathbf{exn} \ l \ \rangle_f :: \psi, \eta) \triangleleft k, \gamma]_o\ $	

(d) C Reductions – $(\tau, \epsilon, \psi, \omega) \xrightarrow{c} \mathfrak{C}$.

ADMINO	$(\tau, \epsilon, (\psi, \eta) \triangleleft k, \gamma) \xrightarrow{o} \ \tau', \epsilon', [(\psi', \eta) \triangleleft k, \gamma]_o\ $	if $(\tau, \epsilon, \psi) \rightsquigarrow (\tau', \epsilon', \psi')$
CALLO	$(v, _, \langle \langle \lambda^o x. e, \epsilon \rangle \rangle_f :: \psi, \eta) \triangleleft k, \gamma) \xrightarrow{o} \ e, \epsilon[x \mapsto v], [(\psi, \eta) \triangleleft k, \gamma]_o\ $	
EXTCALL	$(v, _, \langle \langle \lambda^c x. e, \epsilon \rangle \rangle_f :: \psi, \eta) \triangleleft k, \gamma) \xrightarrow{o} \ e, \epsilon[x \mapsto v], [[], [(\psi, \eta) \triangleleft k, \gamma]_o]_c\ $	
RETTOC	$(v, _, [([], (h, \emptyset)), \gamma) \xrightarrow{o} \ v, \epsilon, \gamma\ $	if $h = \{\mathbf{return} \ x \mapsto x\}$
RETFIB	$(v, _, [([], (h, \epsilon)) \triangleleft k, \gamma) \xrightarrow{o} \ e, \epsilon[x \mapsto v], [k, \gamma]_o\ $	if $\{\mathbf{return} \ x \mapsto e\} \in h$ and $k \neq []$
HANDLE	$(\mathbf{match} \ e \ \mathbf{with} \ h, \epsilon, k, \gamma) \xrightarrow{o} \ e, \epsilon, [([], (h, \epsilon)) \triangleleft k, \gamma]_o\ $	
EXNHN	$(v, _, \langle \langle \mathbf{exn} \ l \ \rangle_f :: _, (h, \epsilon) \rangle \triangleleft k', \gamma) \xrightarrow{o} \ e, \epsilon[x \mapsto v], [k', \gamma]_o\ $	if $\{\mathbf{exception} \ l \ x \mapsto e\} \in h$
EXNFWD C	$(v, \epsilon, \langle \langle \mathbf{exn} \ l \ \rangle_f :: _, (h, _) \rangle, [\psi', \omega]_c) \xrightarrow{o} \ v, \epsilon, [\langle \mathbf{exn} \ l \ \rangle_f :: \psi', \omega]_c\ $	if $\{\mathbf{exception} \ l \ _ \mapsto _ \} \notin h$
EXNFWD FIB	$(v, \epsilon, \langle \langle \mathbf{exn} \ l \ \rangle_f :: _, (h, _) \rangle \triangleleft (\psi', \eta') \triangleleft k', \gamma) \xrightarrow{o} \ v, \epsilon, [(\langle \mathbf{exn} \ l \ \rangle_f :: \psi', \eta') \triangleleft k', \gamma]_o\ $	if $\{\mathbf{exception} \ l \ _ \mapsto _ \} \notin h$
EFFHN	$(v, _, \langle \langle \mathbf{eff} \ l \ k \ \rangle_f :: \psi, (h, \epsilon) \rangle \triangleleft k', \gamma) \xrightarrow{o} \ e, \epsilon[r \mapsto k''] [x \mapsto v], [k', \gamma]_o\ $	if $\{\mathbf{effect} \ l \ x \mapsto e\} \in h$ and $k'' = k @ [(\psi, (h, \epsilon))]$
EFFFWD	$(v, \epsilon', \langle \langle \mathbf{eff} \ l \ k \ \rangle_f :: \psi, (h, \epsilon) \rangle \triangleleft (\psi', \eta') \triangleleft k', \gamma) \xrightarrow{o} \ v, \epsilon', [(\langle \mathbf{eff} \ l \ k'' \ \rangle_f :: \psi', \eta') \triangleleft k', \gamma]_o\ $	if $\{\mathbf{effect} \ l \ _ \mapsto _ \} \notin h$ and $k'' = k @ [(\psi, (h, \epsilon))]$
EFFUNHN	$(v, _, \langle \langle \mathbf{eff} \ l \ k \ \rangle_f :: \psi, (h, \epsilon) \rangle, \gamma) \xrightarrow{o} \ e, \emptyset, [k @ [(\psi, (h, \epsilon))], \gamma]_o\ $	if $\{\mathbf{effect} \ l \ _ \mapsto _ \} \notin h$ and $e = \mathbf{raise} \ \text{Unhandled } 0$
RESUME	$(v, _, \langle \langle k \ \rangle_f :: \langle \langle \lambda^o x. e, \epsilon \rangle \rangle_f :: \psi, \eta) \triangleleft k', \gamma) \xrightarrow{o} \ e, \epsilon[x \mapsto v], [k @ [(\psi, \eta) \triangleleft k'], \gamma]_o\ $	

(e) OCaml Reductions – $(\tau, \epsilon, k, \gamma) \xrightarrow{o} \mathfrak{C}$.

Figure 2. Operational semantics of Multicore OCaml effect handlers.

e) with label l , parameter x and body e , and zero or more effect cases (**effect** $l\ x\ k \mapsto e$) with label l , parameter x , continuation k and body e .

The operational semantics is an extension of the CEK machine semantics [23] for effect handlers, following the abstract machine semantics of Hillerstrom et al. [29]. The key difference from Hillerstrom et al. is that our stacks are composed of alternating sequence of OCaml and C stack segments. The program state is captured as configuration $\mathcal{C} := \|\tau, \epsilon, \sigma\|$ with the current term τ under evaluation, its environment ϵ and the current stack σ . The term is either an expression e or a value v . The values are integer constants n , continuations k , closures $(\lambda x. e, \epsilon)$, effects being performed (**eff** $l\ k$) and exceptions being raised (**exn** l). The environment is a map from variables to values.

The stack σ is either a C stack (γ) or an OCaml stack (ω). The C stack $[\psi, \omega]_c$ consists of a list of frames ψ , and the OCaml stack ω under it. The OCaml stack is either empty \bullet or non-empty $[k, \gamma]_o$ with the current continuation k and the C stack γ under it. Thus, the program stack is an alternating sequence of C and OCaml stacks terminating with an empty OCaml stack \bullet . The frame list ψ is composed of individual frames r , which is one of an argument frame $\langle e \rangle_a$ with the expression e at the argument position of an application with its environment ϵ , a function frame $\langle v \rangle_f$ with the value v at the function position of an application, and frames for evaluating the arguments of an arithmetic expression.

A continuation k is either empty $[]$ or a non-empty list of *fibers*. A fiber $\varphi := (\psi, \eta)$ is a list of frames ψ and a handler closure $\eta := (h, \epsilon)$, which is a pair of handler h and its environment ϵ . We use the infix operator $@$ for appending two lists.

4.2.2 Top-level reductions. The initial configuration for an expression e is $\|e, \emptyset, [[], \bullet]\|$, where the environment and the stack are empty. The top-level reductions (Figure 2b) can be performed by either by taking a C step \xrightarrow{c} or an OCaml step \xrightarrow{o} .

4.2.3 C reductions. We can take a C step (Figure 2d) by taking an administrative reduction step \rightsquigarrow . The administrative reductions are common to both C and OCaml. The rules VAR, ARITH1, ARITH2, APP1, APP2 and APP3 are standard. ARITH3 performs the arithmetic operation on the integers $(\|n_1 \odot n_2\|)$. RAISE pushes a function frame with exception value to indicate that an exception is being raised. Similarly, PERFORM pushes a function frame with an effect value with an empty continuation $[[], (\{\text{return } x \mapsto x\}, \emptyset)]$ with no captured frames and an empty handler with an identity return case alone. We shall return to RESUME1 and RESUME2 later.

Continuing with the rest of the C reduction steps, CALLC captures the behaviour of calling a C function. Since the program is currently executing C, we can perform the call on the

current stack. In case the abstraction is an OCaml abstraction (CALLBACK), we create an OCaml stack with the C stack as its tail, with the current continuation being empty. This captures the behaviour of calling back into OCaml from C. RETTOO returns a value to the enclosing OCaml stack. EXNFWDO forwards a raised exception to the enclosing OCaml stack, unwinding the rest of the frames. This captures the semantics of raising OCaml exceptions from C.

4.2.4 OCaml reductions. In OCaml (Figure 2e), reductions always occur on the top-most fiber in the current stack. ADMINO performs administrative reductions. CALLO evaluates an OCaml function on the current stack. EXTCALL captures the behaviour of external calls, which are evaluated on an empty C stack with the current OCaml stack as its tail. RETTOC returns a value to the enclosing C stack. In this case, we have exactly one fiber on the stack, and this was created in the rule CALLBACK, whose handler has identity return case alone and the environment is empty. RETFIB returns the value from a fiber to the previous one, evaluating the body of the return case.

The rule HANDLE installs a handler by pushing a fiber with no frames and the given handler. The rule EXNHN handles an exception, if the current handler has a matching exception case, unwinding the current fiber. The rule EXNFWDC forwards the exception to C. Here, there is exactly one fiber on the current stack, and the handler does not have a matching exception case, which we know is the case (see CALLBACK rule). The rule EXNFWDFIB forwards the exception to the next fiber if the current handler does not handle it.

The rule EFFHN captures the handling of effects when the current handler has a matching effect case. We evaluate the body of the matching case, and bind the continuation parameter τ to the captured continuation k'' . Observe that the captured continuation k'' includes the current handler. Intuitively, the handler wraps around captured continuation. This gives Multicore OCaml effect handlers deep handler semantics [29]. EFFFWD forwards the effect to the outer fiber, and extends the captured continuation k'' in the process. Recall that we do not capture C frames as part of a continuation. To this end, EFFUNHN models unhandled effect. If the effect bubbles up to the top fiber — which we know does not have an effect case (see CALLBACK rule) — we raise UNHANDLED exception at the point where the corresponding effect was performed. This is achieved by appending the captured continuation to the front of the current continuation.

Observe that **continue** and **discontinue** are not part of the expressions. They are encoded as follows: **continue** $k\ e = (k\ (\lambda^o x. x))\ e$ and **discontinue** $k\ l\ e = (k\ (\lambda^o x. \text{raise } l\ x))\ e$. Intuitively, resuming a continuation in both the cases involves evaluating the appropriate abstraction on top of the continuation. We perform the administrative reductions RESUME1 and RESUME2 to evaluate the arguments to **continue**

and `discontinue`. The rule `RESUME` appends the given continuation to the front of the current continuation, and evaluates the body of the closure.

5 Implementation

We now present the implementation details of effect handlers in Multicore OCaml. We assume an `x86_64` architecture for the remainder of this paper. Although Multicore OCaml currently only supports the `x86_64` backend on Linux, our design does not preclude other CPUs and OSs.

5.1 Exceptions

The implementation follows the operational semantics, but has a few key representational differences. Unlike the operational semantics, handlers with just exception patterns (exception handlers) are implemented differently than effect handlers. As mentioned in §2.2, exceptions are pervasive in OCaml and are so cheap that they are used for local control flow. Hence, we retain the linked exception handler frame implementation of stock OCaml exception handlers in Multicore OCaml for performance backwards compatibility. This differs from other research languages with effect handlers [4, 13, 29], which implement exceptions using effects (by ignoring the continuation argument in the handler).

5.2 Heap-allocated fibers

In the operational semantics, the continuations may be resumed more than once. Captured continuations are copies of the original fibers and resuming the continuation copies the fibers and leaves the continuation as it is. In our use cases, continuations will be resumed at-most once, and copying fibers is unnecessary and inefficient. Instead, Multicore OCaml optimises fibers for one-shot continuations. Fibers are allocated on the C heap using `malloc` and are freed when the handled computation returns with a value or an exception. Similar to Farvardin et al. [22], we use a *stack cache* of recently freed stacks in order to speed up allocation.

Figure 3a shows the layout of a fiber in Multicore OCaml. At the bottom of the stack, we have the `handler_info`, which contains the pointer to the parent fiber, and the closures for the value, exception and effect cases. The closures are created by the translation described in §4.1; Multicore OCaml supports exception patterns in addition to effect patterns in the same handler. This is followed by a context block needed for DWARF and GC bookkeeping with callbacks. Then, there is a top-level exception handler frame that forwards exceptions to the parent fiber. When the exceptions are caught by this handler, the control switches to the parent stack, and the exception handler closure `clos_hexn` is invoked. This is followed by the `pc` of the code that returns values to the parent fiber. This stack is laid out such that when the handled computation returns, the control switches to the parent fiber and the value handler `clos_hval` is invoked.

Next, we have the variable-sized area for the OCaml frames. In order to keep fibers small, this area is initially 16 words in length. When the stack pointer `rsp` becomes less than the *stack threshold* (maintained in the `Cam1_state` table), the stack is said to have overflowed. On stack overflow, we copy the whole fiber to a new area with double the size. In Multicore OCaml, we introduce stack overflow checks into the function prologue of OCaml functions. These stack overflows are rare and so the overflow checks will be correctly predicted by the CPU branch predictor.

In our evaluation of real world OCaml programs (§6), we observed that most function calls are to leaf functions with small frame sizes. Can we eliminate the stack overflow checks for these functions? To this end, we introduce a small, fixed-sized red zone at the top of the stack. The compiler elides the stack overflow check for leaf functions whose frame size is less than the size of the red zone. The default size of the red zone in Multicore OCaml is 16 words.

Finally, we have the saved exception pointer, which points to the top-most exception frame, and the saved stack pointer, which points to the top of the stack. Switching between fibers only involves saving the exception and the stack pointer of the current stack and loading the same on the target stack. Since OCaml does not generate pointers into the stack, the two `fiber_info` fields are the only ones that need to be updated when fibers are moved.

5.3 External calls and callbacks

Since C functions do not have stack overflow checks, we have to execute the external calls in the system stack. Calling a C function from OCaml involves saving the stack pointer in the current fiber, saving the allocation pointer value in `r15` in the `Cam1_state`, updating `rsp` to the top of system stack (maintained in `Cam1_state`), and calling the C function. The actions are reversed when returning from the external call. For C functions that take arguments on the stack, the arguments must be copied to the C stack from the OCaml stack.

When we first enter OCaml from C, a new fiber is allocated for the main OCaml stack. Since callbacks may be frequent in OCaml programs that use finalisers, we run the callbacks on the same fiber as the current one. For example, the layout of the Multicore OCaml stack at `cam1_raise_exn` in the meander example from §2 is shown in Figure 3b. The functions `cam1_call_c` and `cam1_call_ocaml` switch the stacks, and hence are shown in both the system stack and the fiber. Since we are reusing the fiber for the callback, care must be taken to save and restore the `handler_info` before calling and after returning from `c_to_ocaml` function, respectively. Thanks to the fiber representation, external calls and callbacks remain competitive with stock OCaml.

5.4 Effect handlers

In order to illustrate the implementation of handling effects, consider the example presented in Figure 3c. The layout of the program state as the program executes is captured in

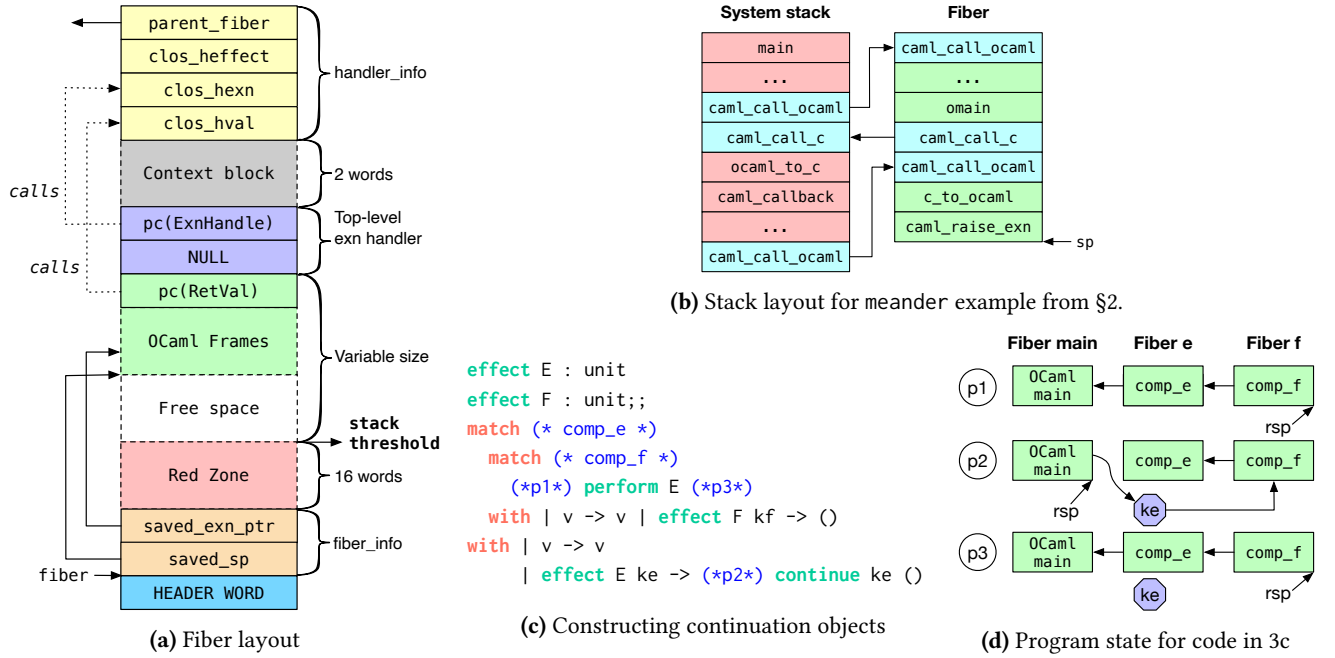


Figure 3. Layout of Multicore OCaml effect handlers.

Figure 3d. The code performs effect E which is handled in the outer-most handler, and is immediately resumed. The arrows between the fibers are parent pointers. At position p_1 , rsp is at the top of the fiber f .

When the effect E is performed, we allocate a continuation object ke in the OCaml heap that points to the current fiber f , set fiber f 's parent pointer to $NULL$, and evaluate the continuation closure $clos_heffect$ on the parent fiber e with the effect E and the continuation ke as arguments. Since the first handler does not handle effect E , the effect is *reperformed* (§4.1) by appending the fiber e to the tail of continuation ke , set fiber e 's parent pointer to $NULL$, and evaluate the current continuation closure on the parent fiber $main$ with E and ke as arguments, which handles E (position p_2). Thus, *continuations are captured without copying frames*. Since every handler closure is evaluated until a matching one is found, the time taken to handle an effect is linear in the number of handlers. We observed that the handler stack is shallow in real programs.

When the continuation is resumed, we overwrite the value of ke to $NULL$ to enforce at-most once semantics. Resuming a continuation involves traversing the linked-list of fibers and making the last fiber point to the current fiber. Just as in the operational semantics, the implementation invokes the appropriate closure to either *continue* or *discontinue* the continuation (position p_3). We perform tail-call optimisation so that resumptions at tail positions do not build up stack.

5.5 Stack unwinding

The challenge with DWARF stack unwinding is to make it aware of the non-contiguous stacks. While the complete details of DWARF stack unwinding is beyond the scope of

the paper, it is beneficial to know how DWARF unwind tables are constructed in order to appreciate our solution. We refer the interested reader to Bastian et al. [3] for a good overview of DWARF stack unwinding.

Logically, DWARF call-frame information maintains a large table which records for every machine instruction where the return address and callee-saved registers are stored. To avoid reifying this large table, DWARF directives represent the table using a compact bytecode representation. This representation is interpreted on demand for building the stack unwind table. Within each function, DWARF maintains a *canonical frame address* (CFA) and is traditionally the stack pointer before entering this function. Hence, on x86-64, where the return address is pushed on the stack on call, the return address is at $CFA - 8$.

The goal is to encode the CFA using DWARF directives when stacks are switched. Recall that stack switching occurs in effect handlers, external calls and callbacks. When starting to run on a new fiber, we insert DWARF bytecode to follow the parent_fiber pointer and dereference the saved_sp to get the CFA ($saved_sp + 8$). During callbacks into OCaml, we save the current system stack pointer in the context block in Figure 3a to identify the CFA in the C stack. DWARF unwinding for external calls is implemented by following a link to the current OCaml stack pointer. With these changes, we get the same backtrace for the meander program from §2.3, modulo runtime system functions due to effect handlers.

We have utilized the DWARF verification tool from [3] to flag possible errors in the hand written DWARF directives where we have altered the OCaml assembler to implement

effect handlers. In three places additional manual verification was required; two were where the stack is intentionally switched between fibers and one was where the control returns to the parent fiber which happened to be resized due to effect handler invocation, which the tool could not track.

Despite the correct DWARF unwind information, using DWARF to record call stack information in `perf` only captures the call stack of the current fiber in Multicore OCaml. Since stack unwinding using DWARF is slow due to bytecode interpretation overhead, `perf` dumps the (user) call stack when sampled [3]. This only includes the frames from the current fiber. This is a limitation of `perf` and not of our stack layout.

5.6 Garbage collection

Recall that OCaml programs are written with the expectation that function calls return exactly once (§3). Consider the scenario when a continuation is neither continued nor discontinued. Since fibers are `malloced`, and are `freed` when the computation being handled returns a value or raises an exception, not resuming continuations leaks memory. In addition, unresumed continuations may also leak other resource such as open file descriptors.

We consider well-behaved Multicore OCaml programs to resume their continuations *exactly once*. This can be enforced by installing a finaliser on every captured continuation that calls `discontinue` on the continuation during finalisation and ignores the result:

```
Gc.finalise (fun k ->
  try ignore (discontinue k Unwind) with _ -> ()) k
```

However, Multicore OCaml does not enforce this behaviour in order to avoid the costs associated with installing a finaliser every time a continuation is captured.

For stack walking during GC, the only change necessary was to make the runtime aware of fiber layout and how to follow the parent pointers. The stack map generation and scanning each frame for roots remains the same. The challenges and solutions for integrating fibers with the concurrent mark-and-sweep GC of Multicore OCaml have been discussed previously [54].

6 Evaluation

In this section, we evaluate the performance of Multicore OCaml effect handlers against the performance requirements set in §1.4. Multicore OCaml is an extension of the 4.10.0 OCaml compiler with support for shared memory parallelism and effect handlers. These results were obtained on a 2-socket Intel®Xeon®Gold 5120 x86-64 [31] server running Ubuntu 18.04 with 64GB of main memory.

6.1 No effects benchmarks

In this section, we measure the impact of the addition of effect handlers on code that does not use effect handlers. Our macro benchmark suite consists of 54 real OCaml workloads including verification tools (Coq, Cubicle, AltErgo), parsers (menhir,

Table 1. Micro benchmarks without effects. Each entry is the percentage difference for Multicore OCaml over stock OCaml.

	<i>exnval</i>	<i>exnraise</i>	<i>extcall</i>	<i>callback</i>	<i>ack</i>	<i>fib</i>	<i>motzkin</i>	<i>sudan</i>	<i>tak</i>
Time	+0.0	-1.9	+17	+65	+5.3	+2.2	+10	+0.0	+4.2
Instr	+0.0	+0.0	+10	+72	+16	+24	+16	+14	+17

`yojson`), storage engines (`irmin`), utilities (`cpdf`, `decompress`), bioinformatics (`fasta`, `knucleotide`, `revcomp2`, `regexredux2`), numerical analysis (`grammatrix`, `LU_decomposition`) and simulations (`nbody`, `game_of_life`). In addition to Stock (stock) and Multicore OCaml (mc), we also ran the benchmarks on Multicore OCaml with no red zone (MC+RedZone0) in the fibers (all OCaml functions will have a stack overflow check) and a red zone size of 32 words (MC+RedZone32). Recall that the default red zone size in Multicore OCaml is 16 words (§5.2).

The complete results are presented in the supplementary material, which we summarise here. On average (geometric mean of the normalized values against stock as the baseline), the multicore variants were **less than 1% slower** than stock. The outliers (on either ends) were due to the difference in the allocator and the GC between stock and Multicore OCaml.

The biggest impact was the increase in the OCaml text section size (OTSS) due to the stack overflow checks. We define OTSS as the sum of the sizes of all the OCaml text sections in the compiled binary file ignoring the data sections, the debug symbols, the text sections associated with OCaml runtime and other statically linked C libraries. Compared to stock, OTSS is 19% more for MC and MC+RedZone32, and 30% more for MC+RedZone0. The result shows that our 16-word red zone is effective at reducing OTSS compared to having no red zone, whereas the 32-word red zone does not noticeably reduce OTSS further. Further work is required to bring OTSS closer to stock.

We also present micro benchmarks results in Table 1. Since micro benchmarks magnify micro-architectural optimisations, we also report the number of instructions executed (obtained using `perf`) along with time. `exnval` performs 100 million iterations of installing exception handlers and returning with a value. `exnraise` is similar, but raises an exception in each iteration. `extcall` and `callback` perform 100 million external calls and callbacks to identity functions. The other micro benchmarks are highly recursive programs and were taken from [22]. For micro benchmarks, we observed that padding tight loops with a few `nop` instructions, which changes the loop alignment, makes the code up to 15% faster. Hence, the difference in running times under 15% may not be statistically significant.

The results show that exceptions are no more expensive in MC compared to stock. In the other programs, MC executes more instructions due to stack overflow checks. The performance impact on callbacks is more significant than external calls.

Table 2. Micro benchmarks with handlers but no perform. Each entry is the slowdown factor (\times times) over its idiomatic implementation in stock OCaml.

	ack	fib	motzkin	sudan	tak
MC	12.25	12.05	11.44	6.74	8.9
monad	348.69	69.77	39.24	33.29	42.79

For callbacks, since we reuse the current fiber stack, we need to ensure it has enough room for inserting additional frames, while `stock` does not need to do this. Callback performance is less important than external calls, which are far more numerous. As seen previously, these results however do not affect macro benchmark performance.

6.2 No perform benchmarks

Next, we aim to quantify the overhead of setting up and tearing down effect handlers compared to a non-tail function call. To this end, we surround the non-tail calls in the recursive micro benchmarks with an effect handler. These programs do not `perform` effects. We also implemented the same benchmarks using a concurrency monad [12] (`monad`) as a proxy for CPS versions. Recall that the OCaml compiler does not use CPS in its IR. In the `monad` version, we use a fork to invoke the non-tail call and use an `MVar` to collect its result.

The results are presented in Table 2. They show that using effect handlers (concurrency monad) is $10.02\times$ ($67.09\times$) more expensive than the idiomatic implementation using non-tail calls. The concurrency monad suffers due to the heap allocation of continuation frames (which need to be garbage collected), whereas effect handlers benefit from stack allocation of the frames. For example, the number of major collections for the `ack` benchmark is 0 for stock OCaml, 1 for `MC` and 112 for `monad`. Our concurrency monad (and other monadic concurrency libraries such as `Lwt` [43] and `Async` [1]) also have other downsides – exceptions, backtraces, and DWARF unwinding are no longer useful due to the lack of a stack.

We note that a compiler that uses CPS IR will be faster than the concurrency monad implementation due to optimisations to reduce the heap allocation of continuation frames. But Farvardin et al. [22] show that CPS with optimisations is still slower than using the call stack.

6.3 Concurrent benchmarks

Next we look at benchmarks that utilize non-local control flow using effect handlers. First, we quantify the overhead of individual operations in effect handling. Consider the following annotated code:

```
effect E : unit
(* a *) match (* b *) perform E (* d *) with
| v -> (* e *) v
| effect E k -> (* c *) continue k ()
```

The sequence `a-b` involves allocating a new fiber and switching to it. `b-c` is performing the effect and handling it. `c-d` is resuming the continuation. `d-e` is returning from the fiber with a value and freeing the fiber. We measured the time

taken to execute these sequences using `perf` support for cycle-accurate tracing on modern Intel processors. We executed 10 iterations of the code, with 3 warm-up runs. We observed that the sequences `a-b`, `b-c`, `c-d` and `d-e` took **23 ns**, **5 ns**, **11 ns** and **7 ns**, respectively. For reference, the idle memory load latency for the local NUMA node is 93.2 ns as measured using the Intel MLC tool [46].

Generators. Generators allow data structures to be traversed on demand. Many languages including JavaScript and Python provide generators as a built-in primitive. Using effect handlers (`MC`), given any data structure (`'a t`) and its iterator (`val iter: 'a t -> ('a -> unit) -> unit`), we can derive its generator function (`val next : unit -> 'a option`). The code is available in the supplementary material. We evaluate the performance of traversing a complete binary tree of depth 25 using this generator. This involves 2^{26} stack switches in total. For comparison, we implemented a hand-written, selective CPSed [48], defunctionalised [17] version (`cps`) and a concurrency monad (`monad`) version of the generator for the tree. Both `cps` and `monad` versions are specialised to the binary tree with the usual caveats of not using the stack for function calls. We observed that the `cps` version was the fastest, thanks to specialisation and hand optimisation. `MC` version was only $2.76\times$ slower than `cps` while being a generic solution, and the `monad` version was $8.69\times$ slower than `cps`.

Chameneos. Chameneos [34] is a concurrency game aimed at measuring context switching and synchronization overheads. Our implementation uses `MVars` for synchronization. We compare effect handler (`MC`), concurrency monad (`monad`) and `Lwt`, a widely used concurrency programming library for OCaml (`lwt`) versions. We observed that `MC` was the fastest, and `monad` (`lwt`) was $1.67\times$ ($4.29\times$) slower than `MC`.

Web server. As one of our primary goals of adding effect handlers to OCaml is to support direct-style concurrency, we implement a full-fledged HTTP/1.1 web server using effect handlers (`MC`). We use effect handlers to implement a user-level threading library that spawns a thread per request. We use `httpaf` [30] for HTTP handling, and `libev` [41] for the eventloop. We compare our implementation against an `Lwt` version (`lwt`) which also uses `httpaf` and `libev`. Unlike using effect handlers, the `Lwt` version is written in monadic style and does not have the notion of a thread per request. For comparison, we include a Go 1.13 version (`go`) that uses the `net/http` [47] package. As both the OCaml versions are single threaded, the Go benchmark is run with `GOMAXPROCS=1`.

The client workload was generated with `wrk2` [60]. We maintain 1k open connections and perform requests for a static web page at different rates, and record the service rate and latency. The throughput and tail latency graphs are given in Figures 4a and 4b. In all the versions, the throughput plateaus at around 30k requests per second. We measure the tail latencies at 2/3rd of this rate (20k requests per second) to simulate optimal load. We observe that both of the OCaml

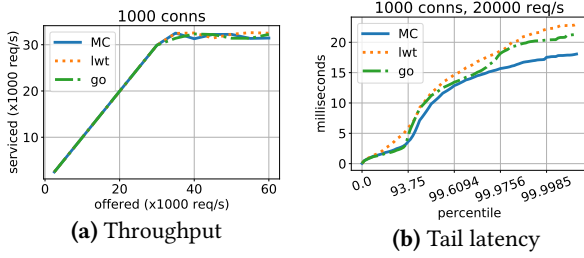


Figure 4. Web server performance.

versions remain competitive with go, and MC performs best in terms of tail latency.

Multicore OCaml supports backtraces for continuations in addition to backtraces of the current stack as in stock OCaml. Using effect handlers in a system such as a web server aides debugging and profiling because it is possible to get a backtrace snapshot of all current requests. Such a feature is available in Go [27], but not in OCaml concurrency libraries such as Async and Lwt which lack the notion of a thread.

7 Related Work

There are several strategies for implementing effect handlers. Eff [4], Helium [7], Frank [13] and the Links server backend [29] use an interpreter similar to our operational semantics to implement effect handlers. Effekt [52], Links JavaScript backend [29] and Koka [38] use type-directed selective CPS translation. These language are equipped with an effect system, which allows compiling pure code in direct-style and effectful code in CPS. Leijen [39] implements effect handlers as a library in C using stack copying. C allows pointers into the stack, so care is taken to ensure that when continuations are resumed, the constituent frames are restored to the same memory addresses as at the time of capture. Kiselyov et al. [37] use an implementation of multi-prompt delimited continuations as an OCaml library [36] to embed the Eff language in OCaml. Indeed, Forster et al. [26] showed that in an untyped setting, effect handlers, monadic reflection and delimited control can macro-express each other.

Multicore OCaml uses the call stack for implementing continuations (as do [36, 39]), but with one-shot continuations. Bruggeman et al. [10] show how to implement one-shot continuations efficiently using segmented stacks in Scheme. Farvardin et al. [22] perform a comprehensive evaluation of various implementation strategies for continuations on modern hardware. Multicore OCaml stacks do not neatly fit the description of one of these implementation strategies – they are best described as using the resize strategy from Farvardin et al. for each of the fibers, which are linked to represent the current stack and the captured continuations.

Kawahara et al. [35] implement one-shot effect handlers using coroutines as a macro-expressible translation, and present an embedding in Lua and Ruby. Lua provides asymmetric coroutines [18] where each coroutine uses its own

stack similar to how each handled computation runs in its own fiber in Multicore OCaml.

Multicore OCaml is not the first language to support stack inspection in the presence of non-local control operators. Chez Scheme supports continuation marks [24] which permit stack inspection as a language feature. This enables implementation of dynamic binding, exceptions, profilers, debuggers, etc, in the presence of first-class continuations. As the authors note, continuation marks can be implemented using effect handlers, but direct support for continuation marks leads to better performance. In this work, we focus on retaining the support for stack inspection through DWARF unwind tables in the presence of effect handlers.

The interaction of non-local control flow and resources has been studied extensively. Scheme uses `dynamic-wind` [51], which is a generalisation of Common Lisp `unwind-protect` [56], which ensures de-allocation (and re-allocation) of resources every time the non-local control leaves and enters back into a context. `dynamic-wind` is not quite sufficient for resource clean up [21, 53] as returning exits must be distinguished from non-returning ones, and resource cleanup is only necessary during non-returning exits.

Multicore OCaml builds on the existing defensive coding practices against exceptions to clean up resources on non-returning exits. Our continuations are assumed to be resumed exactly once using `continue` or `discontinue`. As a result, when a computation performs an effect, we are guaranteed that the control will return. For the non-returning cases (value and exceptional return), the code already handles resource cleanup. OCaml does not have a `try/finally` construct commonly used for resource cleanup in many programming languages. Alternative standard libraries such as Base [2] and Core [15] provide mechanisms analogous to `unwind-protect`, which are implemented using exceptions. Thus, the linear use of continuations enabled by the `discontinue` primitive ensures backwards compatibility of legacy OCaml systems code under non-local control flow introduced by effect handlers. Leijen [40] explicitly extends effect handlers with `initially` and `finally` clauses in Koka for resource safety. Dolan et al. [19] describe the interaction of effect handlers and asynchronous exceptions. This is orthogonal to the contributions of this paper. Our focus is the compiler and runtime system support for implementing effect handlers.

8 Conclusions

Our design for effect handlers in OCaml walks the tightrope of maintaining compatibility (for profiling, debuggers and negligible overheads for existing programs), while unlocking the full power of non-local control flow constructs. Our evaluation shows that we have achieved our goal: we retain compatibility with a surprisingly low performance overhead for sequential code that preserves the spirit of “fast exceptions” that has always characterised OCaml programming. We believe that the introduction of effect handlers into OCaml

implemented using lightweight fibers, along with a parallel runtime [54], will open OCaml to highly scalable concurrent and task-parallel applications with minimal hit to sequential performance.

References

- [1] Async 2020. *Typeful concurrent programming*. <https://opensource.janestreet.com/async/>
- [2] Base.Exn.protect 2020. *Unwind-protect in JaneStreet Base library*. <https://ocaml.janestreet.com/ocaml-core/v0.13/doc/base/Base/Exn/index.html#val-protectx>
- [3] Théophile Bastian, Stephen Kell, and Francesco Zappa Nardelli. 2019. Reliable and Fast DWARF-Based Stack Unwinding. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 146 (Oct. 2019), 24 pages. <https://doi.org/10.1145/3360572>
- [4] Andrej Bauer and Matija Pretnar. 2015. Programming with algebraic effects and handlers. *Journal of Logical and Algebraic Methods in Programming* 84, 1 (2015), 108–123. <https://doi.org/10.1016/j.jlamp.2014.02.001> Special Issue: The 23rd Nordic Workshop on Programming Theory (NWPT 2011) Special Issue: Domains X, International workshop on Domain Theory and applications, Swansea, 5-7 September, 2011.
- [5] Karthikeyan Bhargavan, Barry Bond, Antoine Delignat-Lavaud, Cédric Fournet, Chris Hawblitzel, Catalin Hritcu, Samin Ishtiaq, Markulf Kohlweiss, Rustan Leino, Jay Lorch, Kenji Maillard, Jianyang Pang, Bryan Parno, Jonathan Protzenko, Tahina Ramananandro, Ashay Rane, Aseem Rastogi, Nikhil Swamy, Laure Thompson, Peng Wang, Santiago Zanella-Béguélin, and Jean-Karim Zinzindohoué. 2017. Everest: Towards a Verified, Drop-in Replacement of HTTPS. In *2nd Summit on Advances in Programming Languages*. <http://drops.dagstuhl.de/opus/volltexte/2017/7119/pdf/LIPIcs-SNAPL-2017-1.pdf>
- [6] Dariusz Biernacki, Maciej Piróg, Piotr Polesiuk, and Filip Sieczkowski. 2019. Abstracting Algebraic Effects. *Proc. ACM Program. Lang.* 3, POPL, Article 6 (Jan. 2019), 28 pages. <https://doi.org/10.1145/3290319>
- [7] 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, Article 48 (Dec. 2020), 29 pages. <https://doi.org/10.1145/3371116>
- [8] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horschall, and Noah D. Goodman. 2018. Pyro: Deep Universal Probabilistic Programming. [arXiv:1810.09538 \[cs.LG\]](https://arxiv.org/abs/1810.09538)
- [9] Robert D. Blumofe, Christopher F. Joerg, Bradley C. Kuszmaul, Charles E. Leiserson, Keith H. Randall, and Yuli Zhou. 1995. Cilk: An Efficient Multithreaded Runtime System. In *Proceedings of the Fifth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming* (Santa Barbara, California, USA) (PPOPP '95). Association for Computing Machinery, New York, NY, USA, 207–216. <https://doi.org/10.1145/209936.209958>
- [10] Carl Bruggeman, Oscar Waddell, and R. Kent Dybvig. 1996. Representing Control in the Presence of One-Shot Continuations. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation* (Philadelphia, Pennsylvania, USA) (PLDI '96). Association for Computing Machinery, New York, NY, USA, 99–107. <https://doi.org/10.1145/231379.231395>
- [11] Philippe Charles, Christian Grothoff, Vijay Saraswat, Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. 2005. X10: An Object-Oriented Approach to Non-Uniform Cluster Computing. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications* (San Diego, CA, USA) (OOPSLA '05). Association for Computing Machinery, New York, NY, USA, 519–538. <https://doi.org/10.1145/1094811.1094852>
- [12] Koen Claessen. 1999. A Poor Man's Concurrency Monad. *J. Funct. Program.* 9, 3 (May 1999), 313–323. <https://doi.org/10.1017/S0956796899003342>
- [13] Lukas Convent, Sam Lindley, Conor McBride, and Craig McLaughlin. 2020. Doo bee doo bee doo. *Journal of Functional Programming* 30 (2020), e9. <https://doi.org/10.1017/S0956796820000039>
- [14] Coq 2020. *The Coq Proof Assistant*. <https://coq.inria.fr/>
- [15] Core.Exn.protect 2020. *Unwind-protect in JaneStreet Core library*. <https://ocaml.janestreet.com/ocaml-core/109.20.00/doc/core/Exn.html>
- [16] Olivier Danvy and Andrzej Filinski. 1990. Abstracting Control. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming* (Nice, France) (LFP '90). Association for Computing Machinery, New York, NY, USA, 151–160. <https://doi.org/10.1145/91556.91622>
- [17] Olivier Danvy and Lasse R. Nielsen. 2001. Defunctionalization at Work. In *Proceedings of the 3rd ACM SIGPLAN International Conference on Principles and Practice of Declarative Programming* (Florence, Italy) (PPDP '01). Association for Computing Machinery, New York, NY, USA, 162–174. <https://doi.org/10.1145/773184.773202>
- [18] Ana Lúcia de Moura, Noemi Rodriguez, and Roberto Ierusalimsky. 2004. Coroutines in Lua. *j-jucs* 10, 7 (jul 2004), 910–925. http://www.jucs.org/jucs_10_7/coroutines_in_lua
- [19] Stephen Dolan, Spiros Eliopoulos, Daniel Hillerström, Anil Madhavapeddy, K. C. Sivaramakrishnan, and Leo White. 2018. Concurrent System Programming with Effect Handlers. In *Trends in Functional Programming*, Meng Wang and Scott Owens (Eds.). Springer International Publishing, Cham, 98–117.
- [20] DWARF 2020. *The DWARF Debugging Standard*. <http://dwarfstd.org/>
- [21] Dynamic Wind 2020. *The dynamic-wind problem*. http://okmij.org/ftp/continuations/against-callcc.html#dynamic_wind
- [22] Kavon Farvardin and John Reppy. 2020. From Folklore to Fact: Comparing Implementations of Stacks and Continuations. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) (PLDI 2020). Association for Computing Machinery, New York, NY, USA, 75–90. <https://doi.org/10.1145/3385412.3385994>
- [23] Matthias Felleisen and Daniel P. Friedman. 1986. *Control Operators, the SECD-Machine, and the Lambda-Calculus*. Technical Report. <https://help.luddy.indiana.edu/techreports/TRNNN.cgi?trnum=TR197>
- [24] Matthew Flatt and R. Kent Dybvig. 2020. Compiler and Runtime Support for Continuation Marks. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) (PLDI 2020). Association for Computing Machinery, New York, NY, USA, 45–58. <https://doi.org/10.1145/3385412.3385981>
- [25] Matthew Fluet, Mike Rainey, John Reppy, and Adam Shaw. 2010. Implicitly Threaded Parallelism in Manticore. *J. Funct. Program.* 20, 5–6 (Nov. 2010), 537–576. <https://doi.org/10.1017/S0956796810000201>
- [26] Yannick Forster, Ohad Kammar, Sam Lindley, and Matija Pretnar. 2019. On the expressive power of user-defined effects: Effect handlers, monadic reflection, delimited control. *J. Funct. Program.* 29 (2019), e15. <https://doi.org/10.1017/S0956796819000121>
- [27] Go PProf 2020. *Profiling a Go Program*. <https://golang.org/pkg/runtime/pprof/#Profile>
- [28] Andreas Haas, Andreas Rossberg, Derek L. Schuff, Ben L. Titzer, Michael Holman, Dan Gohman, Luke Wagner, Alon Zakai, and JF Bastien. 2017. Bringing the Web up to Speed with WebAssembly. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Barcelona, Spain) (PLDI 2017). Association for Computing Machinery, New York, NY, USA, 185–200. <https://doi.org/10.1145/3062341.3062363>
- [29] Daniel Hillerström, Sam Lindley, and Robert Atkey. 2020. Effect handlers via generalised continuations. *Journal of Functional Programming* 30 (2020), e5. <https://doi.org/10.1017/S0956796820000040>
- [30] httpaf 2020. *A high performance, memory efficient, and scalable web server written in OCaml*. <https://github.com/inhabitedtype/httpaf>
- [31] Intel Xeon Gold 5120 2020. *Intel® Xeon® Gold 5120 Processor Specification*. <https://ark.intel.com/content/www/us/en/ark/products/120474/intel-xeon-gold-5120-processor-19-25m-cache-2-20-ghz.html>
- [32] IntelTBB 2020. *Intel Thread Building Blocks*. <https://software.intel.com/content/www/us/en/develop/tools/threading-building-blocks.html>
- [33] S. Iwasaki, A. Amer, K. Taura, S. Seo, and P. Balaji. 2019. BOLT: Optimizing OpenMP Parallel Regions with User-Level Threads. In *2019*

- 28th International Conference on Parallel Architectures and Compilation Techniques (PACT). 29–42. <https://doi.org/10.1109/PACT.2019.00011>
- [34] C. Kaiser and J. . Pradat-Peyre. 2003. Chameneos, a concurrency game for Java, Ada and others. In *ACS/IEEE International Conference on Computer Systems and Applications*, 2003. *Book of Abstracts*. 62–. <https://doi.org/10.1109/AICCSA.2003.1227495>
- [35] Satoru Kawahara and Yuki Yoshi Kameyama. 2020. One-Shot Algebraic Effects as Coroutines. In *Trends in Functional Programming*, Aleksander Byrski and John Hughes (Eds.). Springer International Publishing, Cham, 159–179.
- [36] Oleg Kiselyov. 2010. Delimited Control in OCaml, Abstractly and Concretely: System Description. In *Functional and Logic Programming*, Matthias Blume, Naoki Kobayashi, and Germán Vidal (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 304–320.
- [37] Oleg Kiselyov and KC Sivaramakrishnan. 2018. Eff Directly in OCaml. *Electronic Proceedings in Theoretical Computer Science* 285 (Dec 2018), 23–58. <https://doi.org/10.4204/eptcs.285.2>
- [38] Daan Leijen. 2017. Implementing Algebraic Effects in C. In *Asian Symposium on Programming Languages and Systems*, Bor-Yuh Evan Chang (Ed.). Springer International Publishing, Cham, 339–363.
- [39] Daan Leijen. 2017. Type Directed Compilation of Row-Typed Algebraic Effects. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages* (Paris, France) (POPL 2017). Association for Computing Machinery, New York, NY, USA, 486–499. <https://doi.org/10.1145/3009837.3009872>
- [40] Daan Leijen. 2018. *Algebraic Effect Handlers with Resources and Deep Finalization*. Technical Report MSR-TR-2018-10. 35 pages.
- [41] libev 2020. *A high performance full-featured event loop written in C*. <https://metacpan.org/pod/distribution/EV/libev/ev.pod#NAME>
- [42] Loom 2020. *Fibers, continuations and tail-calls for the JVM*. <https://openjdk.java.net/projects/loom/>
- [43] Lwt 2020. *OCaml promises and concurrent I/O*. <https://github.com/ocsigen/lwt>
- [44] Anil Madhavapeddy, Richard Mortier, Charalampos Rotsos, David Scott, Balraj Singh, Thomas Gazagnaire, Steven Smith, Steven Hand, and Jon Crowcroft. 2013. Unikernels: Library Operating Systems for the Cloud. In *Proceedings of the Eighteenth International Conference on Architectural Support for Programming Languages and Operating Systems* (Houston, Texas, USA) (ASPLOS '13). Association for Computing Machinery, New York, NY, USA, 461–472. <https://doi.org/10.1145/2451116.2451167>
- [45] Laurent Mauborgne. 2004. AstrÉE: Verification of Absence of Runtime Error. In *Building the Information Society*, René Jacquart (Ed.). Springer US, Boston, MA, 385–392.
- [46] MLC 2020. *Intel Memory Latency Checker v3.9*. <https://software.intel.com/content/www/us/en/develop/articles/intel-memory-latency-checker.html>
- [47] net/http 2020. *HTTP client and server implementations in Go*. <https://golang.org/pkg/net/http/>
- [48] Lasse R. Nielsen. 2001. A Selective CPS Transformation. *Electronic Notes in Theoretical Computer Science* 45 (2001), 311 – 331. [https://doi.org/10.1016/S1571-0661\(04\)80969-1](https://doi.org/10.1016/S1571-0661(04)80969-1) MFPS 2001, Seventeenth Conference on the Mathematical Foundations of Programming Semantics.
- [49] OCaml Manual 2020. *Extensible variant types*. <https://caml.inria.fr/pub/docs/manual-ocaml/extensiblevariants.html>
- [50] Gordon Plotkin and Matija Pretnar. 2009. Handlers of Algebraic Effects. In *Programming Languages and Systems*, Giuseppe Castagna (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 80–94.
- [51] 1998. Revised5 Report on the Algorithmic Language Scheme. *Higher-Order and Symbolic Computation* 11, 1 (Aug. 1998), 7–105. <https://doi.org/10.1023/A:1010051815785>
- [52] Philipp Schuster, Jonathan Immanuel Brachthäuser, and Klaus Ostermann. 2020. Compiling Effect Handlers in Capability-Passing Style. *Proc. ACM Program. Lang.* 4, ICFP, Article 93 (Aug. 2020), 28 pages. <https://doi.org/10.1145/3408975>
- [53] Dorai Sitaram. 2003. Unwind-protect in portable Scheme. In *Proceedings of the 4th Workshop on Scheme and Functional Programming (7 Nov. 2003)*, M. Flatt, Ed., no. UUCS-03-023 in Tech. Rep., School of Computing, University of Utah. 48–52.
- [54] KC Sivaramakrishnan, Stephen Dolan, Leo White, Sadiq Jaffer, Tom Kelly, Anmol Sahoo, Sudha Parimala, Atul Dhiman, and Anil Madhavapeddy. 2020. Retrofitting Parallelism onto OCaml. *Proc. ACM Program. Lang.* 4, ICFP, Article 113 (Aug. 2020), 30 pages. <https://doi.org/10.1145/3408995>
- [55] KC Sivaramakrishnan, Lukasz Ziarek, and Suresh Jagannathan. 2014. MultiMLton: A multicore-aware runtime for standard ML. *J. Funct. Program.* 24, 6 (2014), 613–674. <https://doi.org/10.1017/S0956796814000161>
- [56] Guy L. Steele. 1990. *Common LISP: The Language (2nd Ed.)*. Digital Press, USA. <https://www.cs.cmu.edu/Groups/AI/html/cltl/cltl2.html>
- [57] Swift 2020. *Swift Concurrency Roadmap*. <https://forums.swift.org/t/swift-concurrency-roadmap/41611>
- [58] Wasm Effect Handlers 2020. *Typed continuations to model stacks*. <https://github.com/WebAssembly/design/issues/1359>
- [59] Sam Westrick, Rohan Yadav, Matthew Fluet, and Umut A. Acar. 2019. Disentanglement in Nested-Parallel Programs. *Proc. ACM Program. Lang.* 4, POPL, Article 47 (Dec. 2019), 32 pages. <https://doi.org/10.1145/3371115>
- [60] Wrk2 2020. *A constant throughput, correct latency recording variant of wrk*. <https://github.com/giltene/wrk2>

Supplementary Material

No effect macro benchmarks

Figures 1 and 2 present the impact on macro benchmarks that do not use effects using the Multicore OCaml compiler which supports effect handlers. The aim is to quantify the impact of the addition of effect handlers to the language on code that does not use effect handlers. The MC is Multicore OCaml with the default red zone size (16 words). MC+RedZone0 uses red zone size of 0; every OCaml function has a stack overflow check. MC+RedZone32 uses a red zone size of 32 words.

Experimental support for multi-shot continuations

Languages with effect handlers typically allow continuations to be resumed more than once to enable use cases such as backtracking search and probabilistic programming (for example, see Uber’s Pyro language [8]). Multicore OCaml has an experimental feature for explicitly cloning continuations, which makes a copy of each of its constituent fibers.

Multi-shot continuations are difficult to reason about in the presence of linear resources such as file descriptors. Recall that OCaml programs are written with the expectation that function calls return *exactly-once* either normally or exceptionally. With multi-shot continuations, they may return *more than once*, leading to errors such as double-freeing malloced memory. Since our goal is to retrofit effect handlers

onto OCaml, we relegate continuation cloning to the OCaml Obj module, which contains unsafe features.

In addition, multi-shot continuations currently are not safe under compiler optimisations. For example consider the following program:

```
effect E : unit
let foo () = perform E
let bar f =
  let r = ref 0 in
  f ();
  r := !r + 1;
  !r
;;
match bar foo with
| v -> v
| effect E k ->
  let k' = Obj.clone_continuation k in
  let r1 = continue k () in
  let r2 = continue k' () in
  Printf.printf "%d_%d\n" r1 r2;
  0
```

This program is expected to print 1 2. However, OCaml optimizes bar by converting the heap allocation of r to a stack allocation. As a result, each continuation ends up with a copy of r, and the program prints 1 1. Thus, in the presence of multi-shot continuations, compiler optimisations need to be made aware of the effects of a function. We leave proper support for multi-shot continuations to future work.

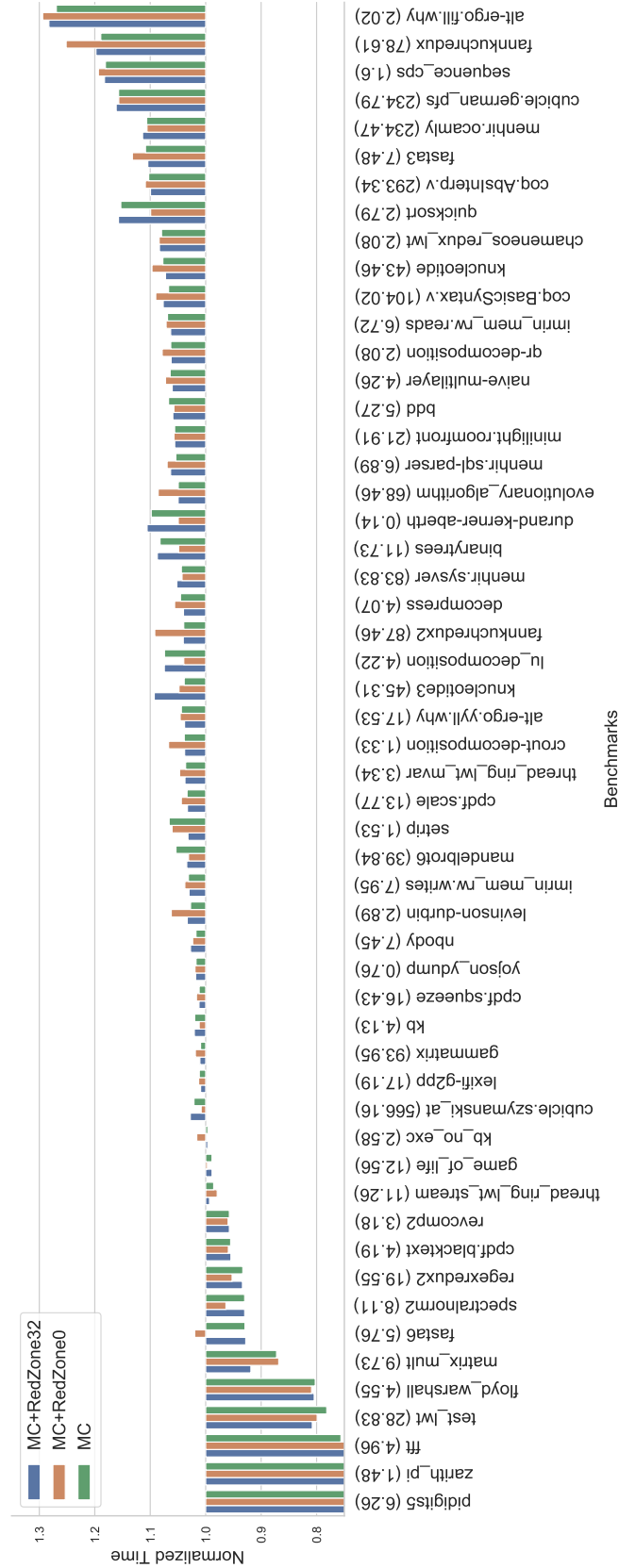


Figure 1. Normalized time of macro benchmarks. Baseline is Stock OCaml, whose running time in seconds is given in parenthesis.

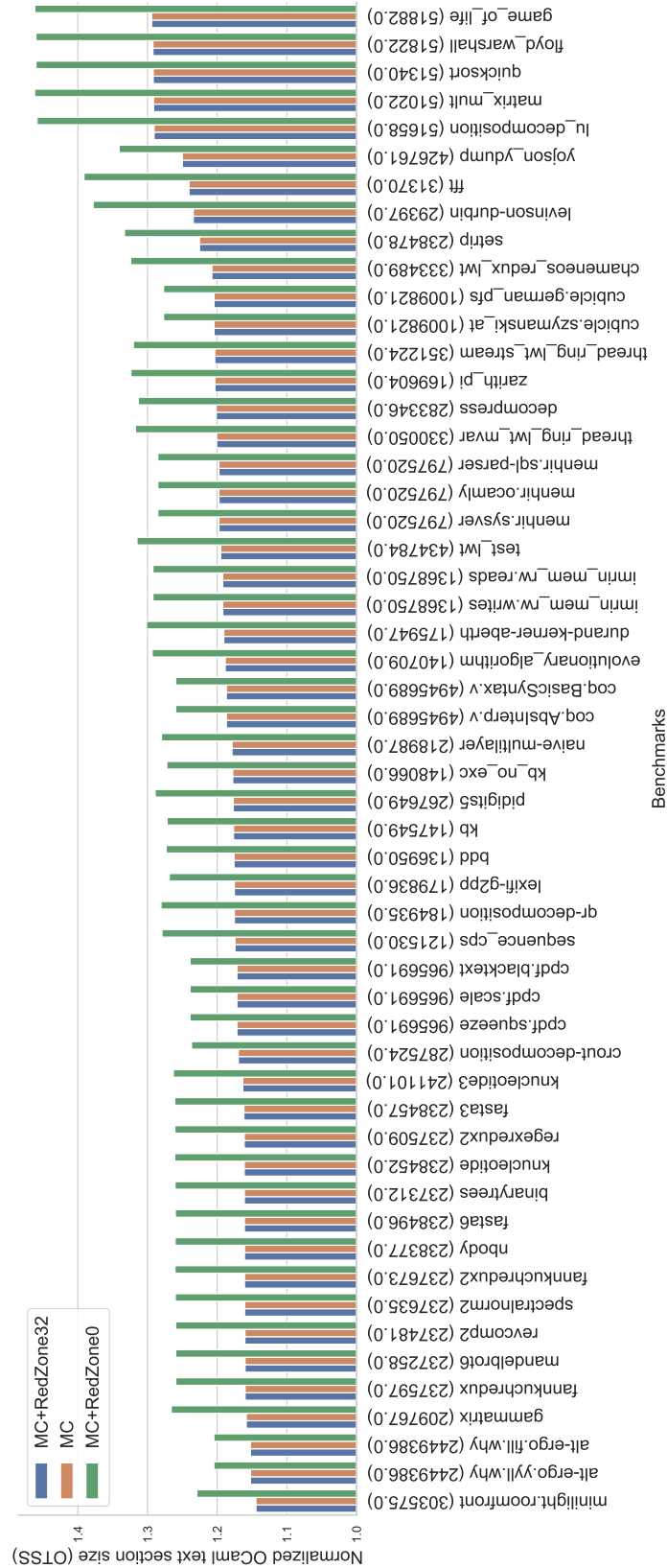


Figure 2. Normalized OCaml text section size (OTSS) of macro benchmarks. Baseline is Stock OCaml, whose size in bytes given in parenthesis.