

Automatic Differentiation via Effects and Handlers in OCaml

Jesse Sigal (University of Edinburgh)

ML 2024, September 6, 2024

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

Automatic differentiation: what and why

- ▶ Derivative based optimization.
- ▶ Automatic differentiation (AD) is a family of algorithms which automatically computes derivatives.
- ▶ AD is only a small constant multiple slower than the original program.
- ▶ Wide variety of implementations and methods.
- ▶ Available methods depend on the language's features.
- ▶ *Idea*: effects and handlers provide a practical basis for AD.

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

Chain rule, two functions

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot f'(x)$$

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

Chain rule, two functions

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot f'(x)$$

Chain rule, three functions

$$\frac{d}{dx}h(g(f(x))) = h'(g(f(x))) \cdot g'(f(x)) \cdot f'(x)$$

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot (f'(x) \cdot 1)$$

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot (f'(x) \cdot 1)$$

```
let (x, dx) = (a, 1) in  
let (y, dy) = (f x, (f' x) * dx) in  
(g y, (g' y) * dy)
```

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in  
let y = f x in  
g y
```

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot (f'(x) \cdot 1)$$

```
let (x, dx) = (a, 1) in  
let (y, dy) = (f x, (f' x) * dx) in  
(g y, (g' y) * dy)
```

$$\frac{d}{dx}g(f(x)) = (1 \cdot g'(f(x))) \cdot f'(x)$$

Automatic differentiation: chain rule two different ways

$g(f(x))$ at $x = a$

```
let x = a in
let y = f x in
g y
```

$$\frac{d}{dx}g(f(x)) = g'(f(x)) \cdot (f'(x) \cdot 1)$$

```
let (x, dx) = (a, 1) in
let (y, dy) = (f x, (f' x) * dx) in
(g y, (g' y) * dy)
```

$$\frac{d}{dx}g(f(x)) = \left(1 \cdot g'(f(x))\right) \cdot f'(x)$$

```
let x = a in
let bz = 1 in
let (z, bx) =
  let y = f x in
  let (z, by) =
    let z = g y in
    (z, bz * (g' y))
  in (z, by * (f' x))
in (z, bx)
```

Automatic differentiation: stateful reverse mode AD

```
let x = a in
let bz = 1 in
let (z, bx) =
  let y = f x in
  let (z, by) =
    let z = g y in
    (z, bz * (g' y))
  in (z, by * (f' x))
in (z, bx)
```

→

```
let (x, bx) = (a, ref 0) in
let (y, by) = (f x, ref 0) in
let (z, bz) = (g y, ref 1) in
by := !by + (!bz * (g' y));
bx := !bx + (!by * (f' x));
(z, !bx)
```

Automatic differentiation: stateful reverse mode AD

```
let x = a in
let bz = 1 in
let (z, bx) =
  let y = f x in
  let (z, by) =
    let z = g y in
    (z, bz * (g' y))
  in (z, by * (f' x))
in (z, bx)
```

\mapsto

```
let (x, bx) = (a, ref 0) in
let (y, by) = (f x, ref 0) in
let (z, bz) = (g y, ref 1) in
by := !by + (!bz * (g' y));
bx := !bx + (!by * (f' x));
(z, !bx)
```

How do we move beyond straight-line programs?

Automatic differentiation: stateful reverse mode AD

```
let x = a in
let bz = 1 in
let (z, bx) =
  let y = f x in
  let (z, by) =
    let z = g y in
    (z, bz * (g' y))
  in (z, by * (f' x))
in (z, bx)
```

\mapsto

```
let (x, bx) = (a, ref 0) in
let (y, by) = (f x, ref 0) in
let (z, bz) = (g y, ref 1) in
by := !by + (!bz * (g' y));
bx := !bx + (!by * (f' x));
(z, !bx)
```

How do we move beyond straight-line programs?

Effects and handlers!

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

Effects and handlers: what and why

- ▶ Structured user-defined side-effects.
- ▶ Like catchable exceptions, but allows continuing from thrown location.
- ▶ Provide abstraction, composition, and reuse.
- ▶ Allows for complicated control flow.
- ▶ In OCaml as of 5.0!
- ▶ Gives a straight-line view of operations.

Effects and handlers: bubbling up

- ▶ We want to “bubble up” the arithmetic operations to get a straight-line program.

Effects and handlers: bubbling up

- ▶ We want to “bubble up” the arithmetic operations to get a straight-line program.
- ▶ Say we have:

```
let x = C[sin 0.5] in D[x]
```

Effects and handlers: bubbling up

- ▶ We want to “bubble up” the arithmetic operations to get a straight-line program.
- ▶ Say we have:

```
let x = C[sin 0.5] in D[x]
```

- ▶ Bubbling up, we get:

```
let y = sin 0.5 in  
let x = C[y] in D[x]
```

Effects and handlers: bubbling up

- ▶ We want to “bubble up” the arithmetic operations to get a straight-line program.
- ▶ Say we have:

```
let x = C[sin 0.5] in D[x]
```

- ▶ Bubbling up, we get:

```
let y = sin 0.5 in  
let x = C[y] in D[x]
```

- ▶ We can also make the second line a first-class function

```
let y = sin 0.5 in  
let k = (fun v -> let x = C[v] in D[x]) in  
k y
```

Effects and handlers: bubbling up

- ▶ We want to “bubble up” the arithmetic operations to get a straight-line program.
- ▶ Say we have:

```
let x = C[sin 0.5] in D[x]
```

- ▶ Bubbling up, we get:

```
let y = sin 0.5 in  
let x = C[y] in D[x]
```

- ▶ We can also make the second line a first-class function

```
let y = sin 0.5 in  
let k = (fun v -> let x = C[v] in D[x]) in  
k y
```

- ▶ Effects and handlers do exactly this!

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

Smooth effect: data types

```
1 type nullary = Const of float
2 type unary = Negate | Sin | Cos | Exp
3 type binary = Plus | Subtract | Times | Divide
4 type arg = L | R
```


Smooth effect: signature

```
6 open Effect
7
8 module type SMOOTH = sig
9   type t
10  type _ Effect.t += Ap0 : nullary -> t Effect.t
11                    | Ap1 : unary * t -> t Effect.t
12                    | Ap2 : binary * t * t -> t Effect.t
13  val c : float -> t
14  val ( ~. ) : t -> t
15  ...
16  val ap0 : nullary -> t
17  val ap1 : unary -> t -> t
18  val ap2 : binary -> t -> t -> t
19  val der1 : unary -> t -> t
20  val der2 : binary -> arg -> t -> t -> t
21 end
```

Smooth effect: base module

```
23 module Smooth (T : sig type t end) : SMOOTH with type t = T.t = struct
24   type t = T.t
25   type _ Effect.t += Ap0 : nullary -> t Effect.t
26                     | Ap1 : unary * t -> t Effect.t
27                     | Ap2 : binary * t * t -> t Effect.t
28   let c x = perform (Ap0 (Const x))
29   let ( ~. ) a = perform (Ap1 (Negate, a))
30   ...
31   let ap0 n = perform (Ap0 n)
32   let ap1 u x = perform (Ap1 (u, x))
33   let ap2 b x y = perform (Ap2 (b, x, y))
34
35   let der1 u x = match u with (*  $\frac{\partial}{\partial x}(u(x))$  *)
36     | Negate -> ~. (c 1.0)      (*  $\frac{\partial}{\partial x}(-x) = -1$  *)
37     | Sin    -> cos_ x         (*  $\frac{\partial}{\partial x}(\sin(x)) = \cos(x)$  *)
38     | Cos    -> ~. (sin_ x)    (*  $\frac{\partial}{\partial x}(\cos(x)) = -\sin(x)$  *)
39     | Exp    -> exp_ x        (*  $\frac{\partial}{\partial x}(e^x) = e^x$  *)
40   ...
```

Smooth effect: base module

```
40 ...
41 let der2 b arg x y = match b with (*  $\frac{\partial}{\partial x_{\text{arg}}}(b(x_L, x_R))$ , for  $x_L = x$ ,  $x_R = y$  *)
42   (*  $\partial/\partial x(x + y) = 1$ ,  $\partial/\partial y(x + y) = 1$  *)
43   | Plus -> (match arg with L -> c 1.0 | R -> c 1.0)
44   (*  $\partial/\partial x(x - y) = 1$ ,  $\partial/\partial y(x - y) = -1$  *)
45   | Subtract -> (match arg with L -> c 1.0 | R -> c (-1.0))
46   (*  $\partial/\partial x(x \cdot y) = y$ ,  $\partial/\partial y(x \cdot y) = x$  *)
47   | Times -> (match arg with L -> y | R -> x)
48   (*  $\partial/\partial x(x/y) = 1/y$ ,  $\partial/\partial y(x/y) = -x/y^2$  *)
49   | Divide ->
50     (match arg with L -> (c 1.0) /. y | R -> (~. x) /. (y *. y))
51 end
```

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

Evaluate: handler

```
1 open Effect.Deep
2 open Float
3 open Smooth
4
5 module Evaluate = struct
6   include Smooth (struct type t = float end)
7
```

Evaluate: handler

```
1 open Effect.Deep
2 open Float
3 open Smooth
4
5 module Evaluate = struct
6   include Smooth (struct type t = float end)
7
8   let (evaluate : ('a, 'a) handler) = {
9     retc = (fun x -> x);
10    exnc = raise;
11    effc = (fun (type x) (eff : x Effect.t) ->
12      match eff with
```

Evaluate: handler

```
1 open Effect.Deep
2 open Float
3 open Smooth
4
5 module Evaluate = struct
6   include Smooth (struct type t = float end)
7
8   let (evaluate : ('a, 'a) handler) = {
9     retc = (fun x -> x);
10    exnc = raise;
11    effc = (fun (type x) (eff : x Effect.t) ->
12      match eff with
13      | Ap0 n -> Some (fun (k : (x, 'a) continuation) ->
14        match n with
15        | Const x -> continue (k : (float, 'a) continuation) x
16      )
17    ...
```

Evaluate: handler

```
16   ...
17   | Ap1 (u, x) -> Some (fun k ->
18       match u with
19         | Negate -> continue k (neg x)
20         | Sin    -> continue k (sin x)
21         | Cos    -> continue k (cos x)
22         | Exp    -> continue k (exp x)
23       )
24   | Ap2 (b, x, y) -> Some (fun k ->
25       match b with
26         | Plus      -> continue k (add x y)
27         | Subtract -> continue k (sub x y)
28         | Times    -> continue k (mul x y)
29         | Divide   -> continue k (div x y)
30       )
31   | _ -> None
32 )
33 }
34 end
```


Evaluate: example

```
1 open Effect.Deep
2 open Evaluate
3
4 let _ =
5   let open Evaluate in
6     let sqr x = x *. x in
7     let res = (match_with : ('c -> 'a) -> 'c -> ('a, 'b) handler -> 'b)
8       (fun (twice, x) -> if twice then sqr (sqr x) else sqr x)
9       (true, 5.0)
10      evaluate
11   in
12   Printf.printf "%f\n" res (* Prints "625.000000"=54 *)
```

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

$$\frac{d}{dx}g(f(x)) = (1 \cdot g'(f(x))) \cdot f'(x)$$

```
let (x, bx) = (a, ref 0) in
let (y, by) = (f x, ref 0) in
let (z, bz) = (g y, ref 1) in
by := !by + ((g' y) * !bz);
bx := !bx + ((f' x) * !by);
(z, !bx)
```

Reverse mode: numeric type

```
1 open Effect.Deep
2 open Smooth
3
4 type 't mpaired = {v : 't; mutable bv : 't}
5
```

Reverse mode: numeric type

```
1 open Effect.Deep
2 open Smooth
3
4 type 't mpaired = {v : 't; mutable bv : 't}
5
6 module Reverse (T : SMOOTH) = struct
7   include Smooth (struct type t = T.t mpaired end)
8
```

Reverse mode: numeric type

```
1 open Effect.Deep
2 open Smooth
3
4 type 't mpaired = {v : 't; mutable bv : 't}
5
6 module Reverse (T : SMOOTH) = struct
7   include Smooth (struct type t = T.t mpaired end)
8
9   let (reverse : (unit, unit) handler) = {
10     retc = (fun x -> x);
11     exnc = raise;
12     effc = (fun (type a) (eff : a Effect.t) ->
13       match eff with
14     ...
```

Reverse mode: handler

```
14  ...
15  | Ap0 n -> Some (fun (k : (a, _) continuation) -> let open T in
16    continue k {v = ap0 n; bv = c 0.0}
17  )
18  | Ap1 (u, x) -> Some (fun k -> let open T in
19    let r = {v = ap1 u x.v; bv = c 0.0} in
20    continue k r;
21    x.bv <- x.bv +. (der1 u x.v *. r.bv)
22  )
23  | Ap2 (b, x, y) -> Some (fun k -> let open T in
24    let r = {v = ap2 b x.v y.v; bv = c 0.0} in
25    continue k r;
26    x.bv <- x.bv +. (der2 b L x.v y.v *. r.bv);
27    y.bv <- y.bv +. (der2 b R x.v y.v *. r.bv)
28  )
29  | _ -> None
30  )
31  }
```

Reverse mode: derivative function

```
31   ...
32   (* grad f x =  $\frac{\partial f(z)}{\partial z}(x)$  *)
33   let grad (f : T.t mpaired -> T.t mpaired) (x : T.t) =
34     let r = {v = x; bv = T.c 0.0} in
35     match_with (fun x -> (f x).bv <- T.c 1.0) r reverse;
36     r.bv
37 end
```


Reverse mode: example

```
1 let _ =
2   let module E = Evaluate in
3   let module R = Reverse(E) in
4   let sqr x = R.(x *. x) in
5   let res = match_with
6     (fun (twice, y) ->
7       R.grad (fun x -> if twice then sqr (sqr x) else sqr x) y
8     )
9     (true, 5.0)
10    E.evaluate
11  in
12  Printf.printf "%f\n" res (* Prints "500.000000" =  $4 \cdot 5^3 = \frac{\partial(x^4)}{\partial x}(5)$  *)
```

Reverse mode: nested example

```
1 let _ =
2   let module E = Evaluate in
3   let module R = Reverse(E) in
4   let module RR = Reverse(R) in
5   let sqr x = RR.(x *. x) in
6   let res = match_with (fun (twice, z) ->
7     R.grad (fun y ->
8       RR.grad (fun x -> if twice then sqr (sqr x) else sqr x) y
9       ) z
10    ) (true, 5.0) E.evaluate
11 in
12 Printf.printf "%f\n" res (* Prints "300.000000" =  $12 \cdot 5^2 = \frac{\partial^2(x^4)}{\partial x^2}(5)$  *)
```

Reverse mode: original implementation

- ▶ We are not the first to implement reverse mode AD with handlers, see [Sivaramakrishnan, 2018].
- ▶ [Sivaramakrishnan, 2018] was inspired by [Wang et al., 2019] who used delimited continuations.
- ▶ We are the first to design a larger system and add tensor valued operations, as well as benchmark.

Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

What to expect

- ▶ [Griewank and Walther, 2008, Sec. 4.4] show that for a composite measure of “work”, reverse mode is $O(1)$ w.r.t. the original program.
- ▶ Work includes
 - ▶ memory fetches and stores,
 - ▶ additions and subtractions,
 - ▶ multiplications, and
 - ▶ non-linear operations.
- ▶ With reasonable assumptions, they prove reverse mode should be $3\times$ to $4\times$ slower.

Microbenchmark: code

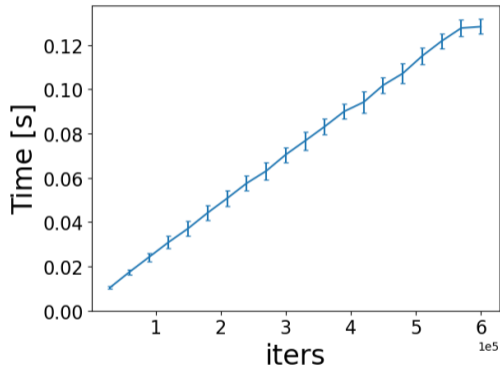
$$\frac{1}{x} = \sum_{n=0}^{\infty} (-1)^n (x-1)^n = \sum_{n=0}^{\infty} a_n$$
$$a_0 = 1, \quad a_n = -(x-1) \cdot a_{n-1}$$

Microbenchmark: code

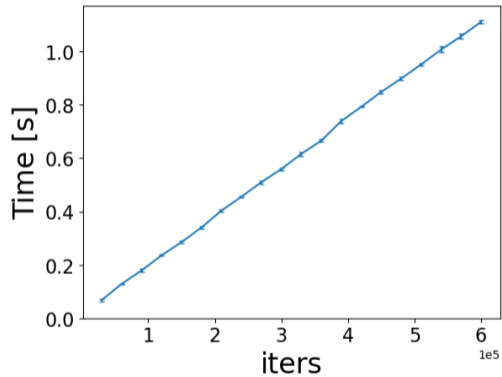
$$\frac{1}{x} = \sum_{n=0}^{\infty} (-1)^n (x-1)^n = \sum_{n=0}^{\infty} a_n$$
$$a_0 = 1, \quad a_n = -(x-1) \cdot a_{n-1}$$

```
1 open Smooth
2
3 module Taylor_Recip_Benchmark (T : SMOOTH) = struct
4   let approx_recip iters x = let open T in
5     let prev = ref (c 1.0) in (* a0 *)
6     let acc = ref (c 1.0) in (* ∑n=00 an *)
7     for _i = 1 to iters do
8       prev := !prev *. (~. (x -. (c 1.0))); (* ai = -(x-1) · ai-1 *)
9       acc := !prev +. !acc (* ∑n=0i an = ai + ∑n=0i-1 an *)
10    done;
11    !acc (* ∑n=0iters an *)
12 end
```

Microbenchmark: results



(a) Evaluation mode



(b) Reverse mode

Microbenchmark: results

Reverse mode is about $8.3\times$ slower

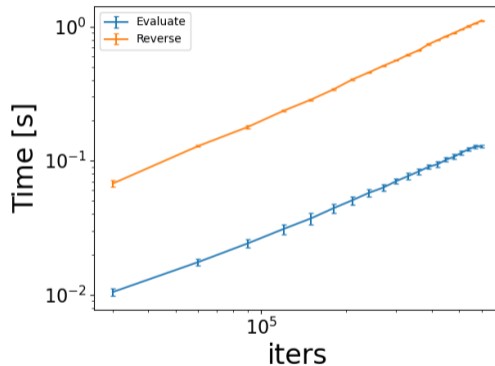


Figure: Reverse and evaluation modes, log-log scale

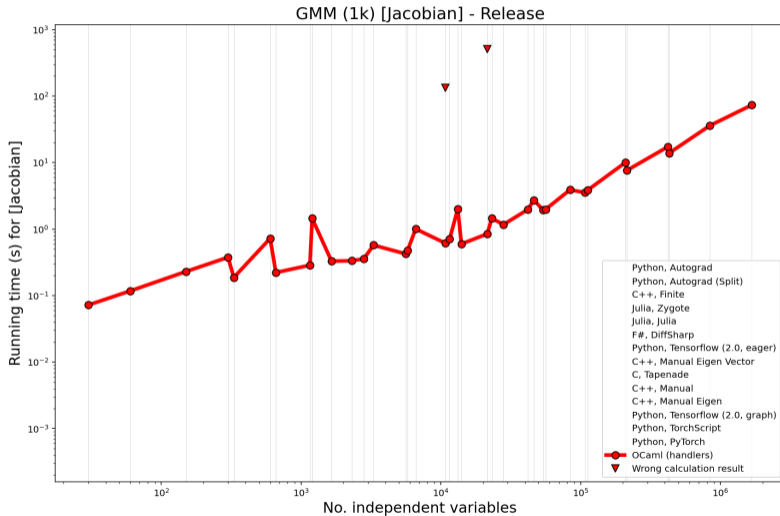
Macrobenchmark

- ▶ Benchmark suite of [Srajer et al., 2018]
 - ▶ reproducible: containers
 - ▶ extensible: test harnesses, modular
 - ▶ realistic: real ML and computer vision functions
- ▶ *Problem*: one effect call per real-valued operation will be inefficient
- ▶ *Solution*: tensor/matrix/ND-array operations
 - ▶ Owl scientific computing library of [Wang et al., 2022]
 - ▶ Extend `Smooth` to 35 operations
 - ▶ Extend `Reverse` to handle new operation types
- ▶ We implement the objective function for Gaussian mixture models
- ▶ 3 different parameters N , K , and D
- ▶ $K \cdot D$ is the total number of input variables, giving our x -axis

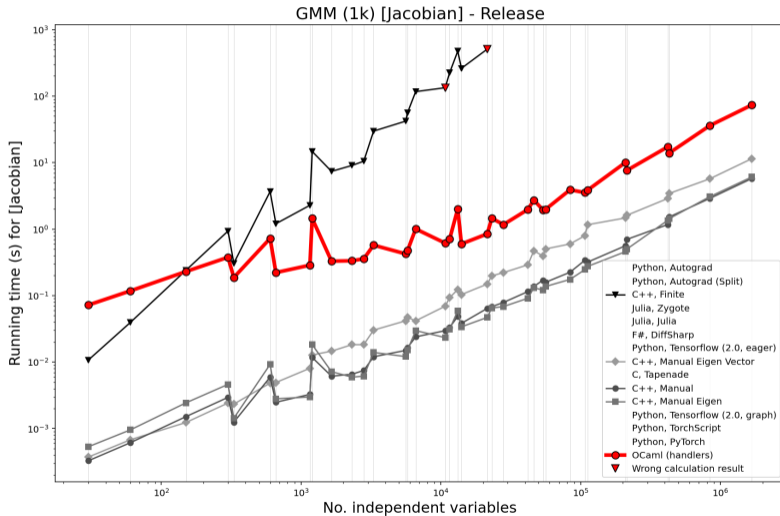
Macrobenchmark: other systems

Language	Tool	Approach
C++	-	Manual (by hand)
C++	-	Finite differences
C	Tapenade	Static
Python	Autograd	Dynamic
Python	TensorFlow 2.0 (eager)	Dynamic
Python	TensorFlow 2.0 (graph)	Static
Python	PyTorch	Dynamic
Python	TorchScript	Static
Julia	ForwardDiff.jl	Dynamic
Julia	Zygote	Static
F#	DiffSharp	Dynamic
OCaml	This work	Dynamic

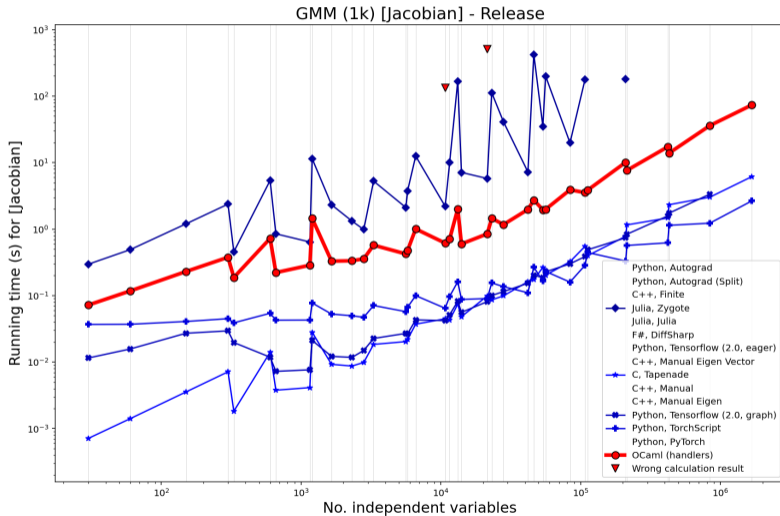
Macrobenchmark: results (1k)



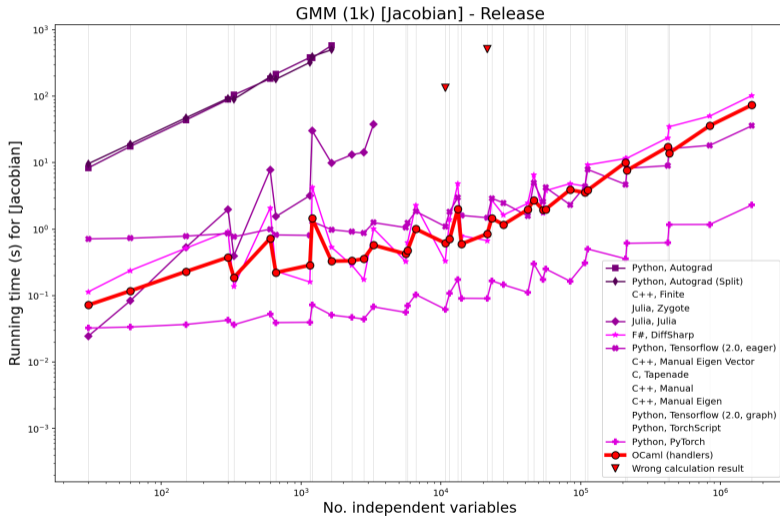
Macrobenchmark: results (1k, manual)



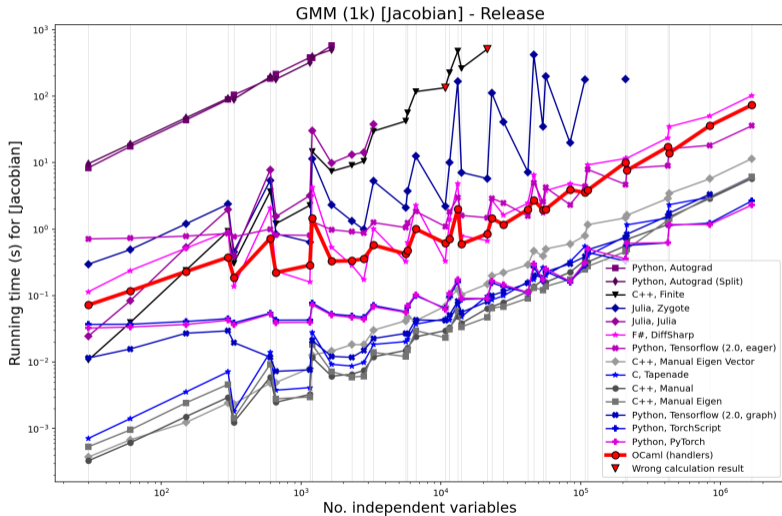
Macrobenchmark: results (1k, static)



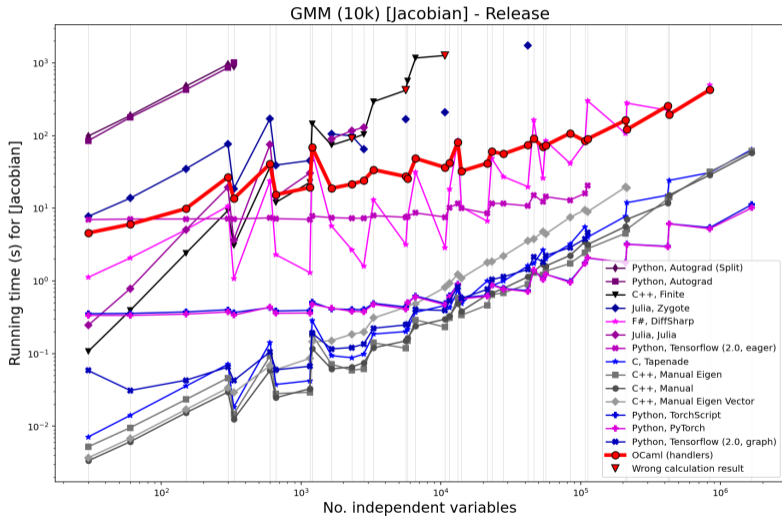
Macrobenchmark: results (1k, dynamic)



Macrobenchmark: results (1k, all)



Macrobenchmark: results (10k)



Overview

Automatic differentiation

Effects and handlers

The `Smooth` effect

Evaluate handler

Reverse mode handler

Benchmarks

Conclusion

More to say on AD and effect handlers

- ▶ Many different modes
 - ▶ Checkpointed reverse mode for time-space tradeoff
 - ▶ Higher-order functions
 - ▶ Hessians
- ▶ Different languages, some with effect type systems
 - ▶ Koka
 - ▶ Frank
 - ▶ Eff
- ▶ Mathematical correctness
 - ▶ Denotational semantics
 - ▶ For forward mode and (simpler) reverse mode

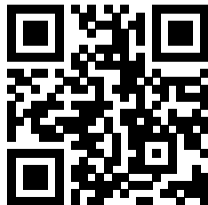
See my thesis:








Conclusion

- ▶ AD (and reverse mode specifically) requires complex control flow.
- ▶ Effect handlers enable a simple implementation which follows the math.
- ▶ With little effort and moving to tensor valued operations, we are competitive among similar tools.
- ▶ **Future work:**
 - ▶ Use Torch bindings in OCaml.
 - ▶ Correctness of reverse mode.
 - ▶ Custom functions (higher-order effects?).

Preprint:



-  Griewank, A. and Walther, A. (2008).
Evaluating Derivatives.
Other Titles in Applied Mathematics. Society for Industrial and Applied Mathematics.
-  Sivaramakrishnan, K. C. (2018).
Reverse-mode Algorithmic differentiation using effect handlers.
-  Srajer, F., Kukelova, Z., and Fitzgibbon, A. (2018).
A benchmark of selected algorithmic differentiation tools on some problems in computer vision and machine learning.
Optimization Methods and Software, 33(4-6):889–906.
Publisher: Taylor & Francis .eprint:
<https://doi.org/10.1080/10556788.2018.1435651>.

-  Wang, F., Zheng, D., Decker, J., Wu, X., Essertel, G. M., and Rompf, T. (2019). Demystifying differentiable programming: shift/reset the penultimate backpropagator.
Proceedings of the ACM on Programming Languages, 3(ICFP):96:1–96:31.
-  Wang, L., Zhao, J., and Mortier, R. (2022).
OCaml Scientific Computing: Functional Programming in Data Science and Artificial Intelligence.
Undergraduate Topics in Computer Science. Springer International Publishing, Cham.