

Zygote

Bridging Machine Learning and Scientific Computing

mike.j.innes@gmail.com



A powerful, high level language with high performance.

Pythonic, mathematical syntax that looks like notation.

Performance consistently within 2x of tuned C code.

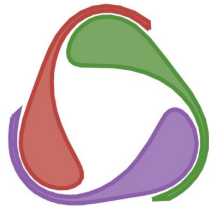
Most of Julia is written in Julia!

```
function mandel(z)
    c = z
    maxiter = 80
    for n = 1:maxiter
        if abs(z) > 2
            return n-1
        end
        z = z^2 + c
    end
    return maxiter
end
```

Scientific Computing



DSGE.jl



DifferentialEquations.jl



Machine Learning

High-level and flexible (Python)

High overhead, focus on tensor operations and manual vectorisation

Relatively simple programs (network architectures)

Mutation support considered advanced/unusual.

Scientific Computing

Low-level and manual (Fortran)

Low overhead, focus on scalar operations

Regularly run over millions of lines of code.

Research on auto-vectorisation, shared memory parallelism, checkpointing etc.



Tapenade



Ruthless pragmatism and scalability.
Output can be highly optimised using
existing optimising compilers.

λ the Ultimate Backpropagator



Elegant recursive formalism, including
nested AD (closure), convenience
(callee-derives) and bags of expressive
power.

```
[(v1.2) pkg> add Zygote
```

```
Resolving package versions...
```

```
Updating `~/julia/environments/v1.2/Project.toml`
```

```
[e88e6eb3] + Zygote v0.3.2
```

```
Updating `~/julia/environments/v1.2/Manifest.toml`
```

```
[1a297f60] + FillArrays v0.6.3
```

```
[7869d1d1] + IRTools v0.2.2
```

```
[e88e6eb3] + Zygote v0.3.2
```

```
[julia> using Zygote
```

```
[julia> function pow(x, n)
```

```
    r = 1
```

```
    while n > 0
```

```
        n -= 1
```

```
        r *= x
```

```
    end
```

```
    return r
```

```
end
```

```
pow (generic function with 1 method)
```

```
[julia> pow(5, 3)
```

```
125
```

```
[julia> gradient(x -> pow(x, 3), 5)
```

```
(75,)
```

```
julia> █
```


User Function	Primal	Adjoint
<pre> function pow(x, n) r = 1 while n > 0 n -= 1 r *= x end return r end </pre>	<pre> 1: (%2, %3) br 2 (%3, 1) 2: (%4, %5) %6 = %4 > 0 br 4 unless %6 br 3 3: %7 = %4 - 1 %8 = %5 * %2 br 2 (%7, %8) 4: return %5 </pre>	<pre> 1: (%1) br 2 (%1, 0) 2: (%2, %4) br 4 unless @6 br 3 3: %10 = %2 * @2 %11 = %2 * @5 %14 = %4 + %11 br 2 (%10, %14) 4: return (%4, 0) </pre>

$\text{pow}(5, 3) == 125$
 $\text{gradient}(\text{pow}, 5, 3) == (75, 0)$

```
[julia> y, back = J(pow, 5, 3);
```

```
[julia> back(1)  
(75, nothing)
```

$$y = f(x_1, x_2, \dots)$$

$$y, \mathcal{B} = \mathcal{J}(f, x_1, x_2, \dots)$$

$$\bar{x}_1, \bar{x}_2, \dots = \mathcal{B}(\bar{y})$$

```
function foo(x)
  a = bar(x)
  b = baz(a)
  return b
end
```



```
function J(::typeof(foo), x)
  a, da = J(bar, x)
  b, db = J(baz, a)
  return b, function(b-)
    ā = db(b-)
    x- = da(ā)
    return x-
  end
end
```

→ ~ j



Documentation: <https://docs.julialang.org>

Type "?" for help, "]? " for Pkg help.

Version 1.2.0-rc1.2 (2019-05-31)

release-1.2/3fcb168ceb (fork: 74 commits, 81 days)

```
[julia> fs = Dict("sin" => sin, "cos" => cos, "tan" => tan);
```

```
[julia> f(x) = fs[readline()](x)
f (generic function with 1 method)
```

```
[julia> f(1)
sin
0.8414709848078965
```

```
[julia> gradient(f, 1)
sin
(0.5403023058681398,)
```

```
julia> █
```

`J(::typeof(sin), x) = sin(x), \bar{y} -> \bar{y} *cos(x)`



`@adjoint sin(x) = sin(x), \bar{y} -> \bar{y} *cos(x)`

Core compiler pass is ~200 lines of code

All semantics added via custom adjoints –
mutation, data structures, checkpointing, etc.

```
nestlevel() = 0
```

```
@adjoint nestlevel() = nestlevel()+1, _ -> nothing
```

```
julia> function f(x)
    println(nestlevel(), " levels of nesting")
    return x
end
```

```
julia> f(1);
0 levels of nesting
```

```
julia> grad(f, 1);
1 levels of nesting
```

```
julia> grad(x -> x*grad(f, x), 1);
2 levels of nesting
```

```
@adjoint hook(f, x) = x, Δ -> (f(Δ),)
```

```
hook(-, x) # reverse the gradient of x
```

```
@adjoint checkpoint(f, x...) =  
    f(x...), Δ -> J(f, x...)[2](Δ)
```

```
@adjoint function forwarddiff(f, x)
```

```
    y, J = forward_jacobian(f, x)
```

```
    y, Δ -> (J'Δ,)
```

```
end
```

```
 julia> hook(f, x) = x
 hook (generic function with 1 method)
```

```
 julia> @adjoint hook(f, x) = x, Δ -> (nothing, f(Δ),)
```

```
 julia> gradient(2, 3) do a, b
     a*b
 end
 (3, 2)
```

```
 julia> gradient(2, 3) do a, b
     hook(-, a) * b
 end
 (-3, 2)
```

```
 julia> gradient(2, 3) do a, b
     hook( $\bar{a}$  -> @show( $\bar{a}$ ), a) * b
 [
     end
  $\bar{a}$  = 3
 (3, 2)
```


Differentiation à la Carte

- Mixed-mode AD (forward, reverse, Taylor series, ...)
- Forward-over-reverse (Hessians)
- Cross-language AD
- Support for Complex and other number types
- Easy custom gradients
- Checkpointing
- Gradient hooks
- Custom types (colours!)
- Hardware backends: CPU, CUDA, TPU, ...
- ~~Deeply nested AD~~ (WIP)

```
196 dense(W, b, σ = identity) =
197     x → σ.(W * x .+ b)
198
199 chain(f ... ) = foldl(∘, reverse(f))
200
201 mlp = chain(
202     dense(randn(5, 10), randn(5), tanh),
203     dense(randn(2, 5), randn(2)))
204
205 x = rand(10)
```

```
207 mlp(x)
```

```
  ▾ Float64[2]
    0.646...
    2.51...
```

Deep learning in 5 lines.

```
211 m̄ = gradient(mlp) do m
212     sum(m(x))
213 end  [ ((f = (W = [-0.9909137325976834  0.11388709497399903 ... -0.7210152885786678  0.99010
214
215 m -= η * m̄ # Gradient descent
```

Data Structures & Mutation

```
julia> using Colors
```

```
julia> a, b = RGB(1, 0, 0), RGB(0, 1, 0)
(RGB{N0f8}(1.0,0.0,0.0), RGB{N0f8}(0.0,1.0,0.0))
```

```
julia> a.r^2
1.0N0f8
```

```
julia> gradient(c -> c.r^2, a)
((r = 2.0f0, g = nothing, b = nothing),)
```

```
julia> colordiff(a, b)
86.60823557376344
```

```
julia> gradient(a -> colordiff(a, b), a)
((r = 0.4590887719632896, g = -9.598786801605689, b = 14.181383399012862),)
```

```
julia> vars = Dict{:r => 0, :n => 0}
Dict{Symbol,Int64} with 2 entries:
  :n => 0
  :r => 0
```

```
julia> function pow(x, n)
    vars[:r] = 1
    vars[:n] = n
    while vars[:n] > 0
        vars[:n] -= 1
        vars[:r] *= x
    end
end
pow (generic function with 1 method)
```

```
julia> pow(5, 3); vars[:r]
125
```

```
julia> gradient(x -> (pow(x, 3); vars[:r]), 5)
(75,)
```

```
julia> vars[:r]
125
```



```
@adjoint function pycall(f, x...; kw...)
  x = map(py, x)
  y = pycall(f, x...; kw...)
  y.detach().numpy(), function ( $\bar{y}$ )
    y.backward(gradient = py( $\bar{y}$ ))
    (nothing, map(x -> x.grad.numpy(), x)...)
end
end
```

Some Bonus Features


```
6 @grad function (a::Real * b::Real)
7   c = a*b
8   function back(Δ)
9     0//0
10  end
11  c, back
12 end |> _forward
```

```
13
14 function pow(x, n)
15   r = one(x)
16   while n > 0
17     r *= x
18     n -= 1
19   end
20   return r
21 end |> pow
```

```
22
23 gradient(pow, 2, 3)
```

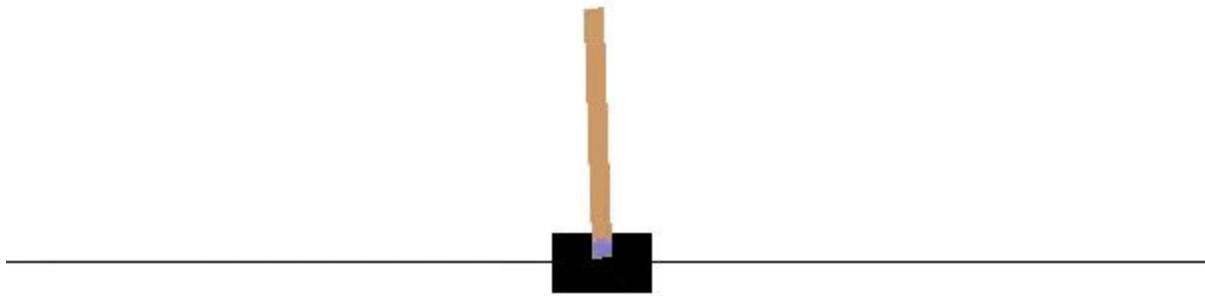
```
24
25 ArgumentError: invalid rational: zero(Int64)//zero(Int64)
26   in top-level scope at base/none
27   in gradient at Zygote/src/compiler/interface.jl:34
28   in at Zygote/src/compiler/interface.jl:28
29   in at Zygote/src/compiler/interface2.jl
30   in pow at test.jl:17
31   in at Zygote/src/lib/lib.jl:33
32   in at test.jl:9
33   in // at base/rational.jl:13
```

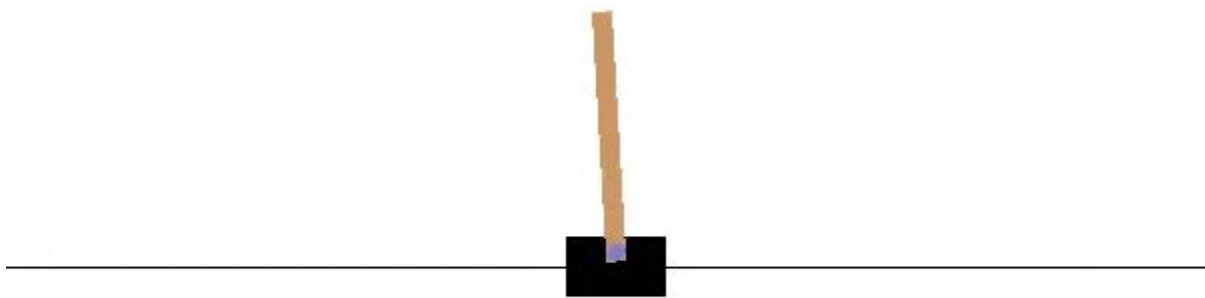
```
8 using Zygote
9
10 function f(x)
11     for i = 1:5
12         x = sin(cos(x))
13     end
14     return x
15 end
16
17 function loop(x, n)
18     r = x/x
19     for i = 1:n
20         r *= f(x)
21     end
22     return sin(cos(r))
23 end
24
25 gradient(loop, 2, 3)
26
27 Zygote.@profile loop(2, 3)
28
29 function logsumexp(x::Array{Float64,1})
```

Future Challenges

- Mutation of values is hard
- Need adjoints to cover the entire standard library
- Compiler improvements
 - More functional-style optimisations
 - Better heuristics for AD-generated code
- Fast code vs. dynamic semantics
- Differentiating Julia's concurrency and parallelism constructs
- Reducing overheads: currently ~50ns per operation
 - Great compared to ML frameworks but far from optimal

Differentiating CartPole





CartPole Controller

CartPole State

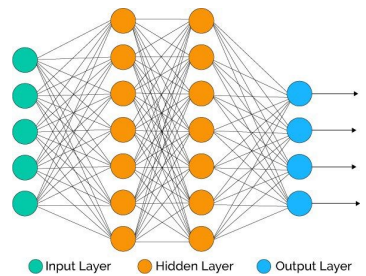
Control Parameters

Loss

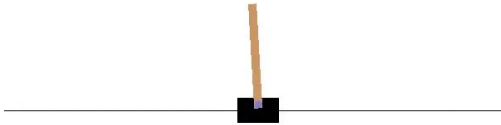
Neural Network

Environment

$angle = -3^\circ$
 $velocity = 0.5^\circ/s$



{left, right}



$angle^2$



Backpropagation

Model-free RL

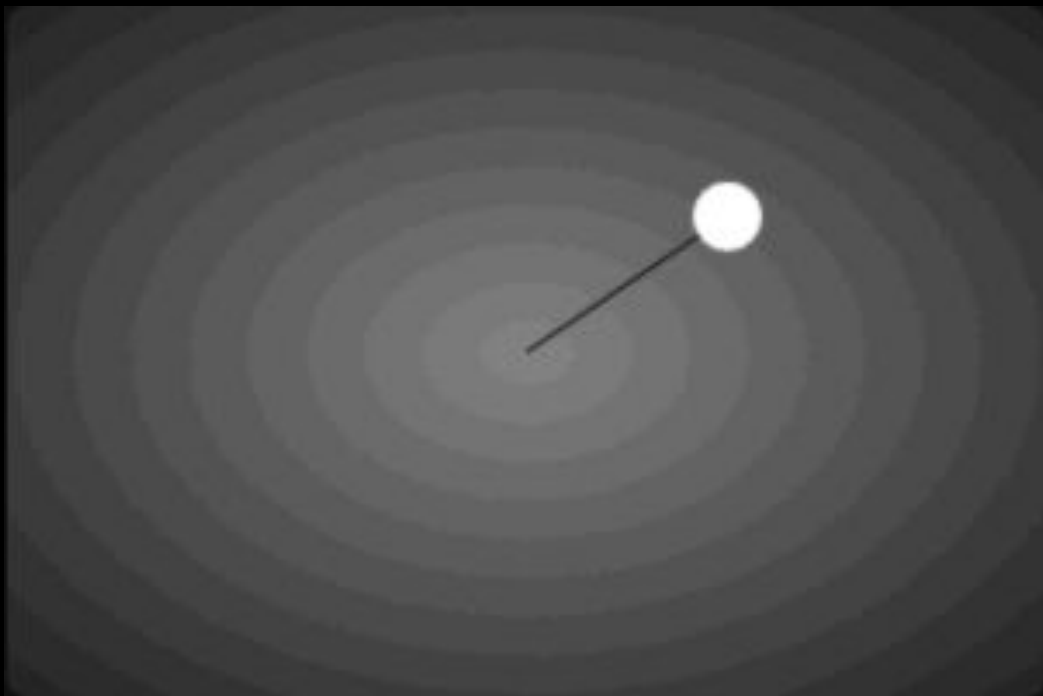
~400

episodes to solve

Differentiable Programming

6

episodes to solve



Pharmacokinetics/Pharmacodynamics

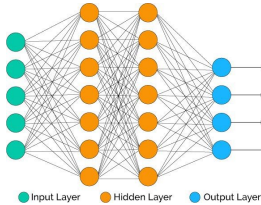
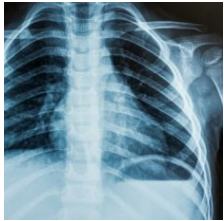
Patient Data
(e.g. X-ray scans)

Simulation parameters
(e.g. drug half-life)

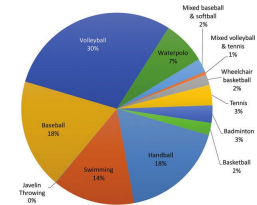
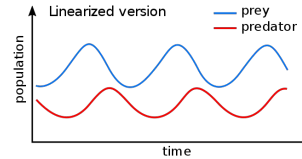
Treatment Outcome

Neural Network

ODE Solver



$$\beta = 0.56$$



Backpropagation

Is Differentiable Programming a Thing?

Bona fide programming paradigm, or yet another rebranding of neural networks?

Programming Paradigms

A collection of design patterns and organising principles, built on a common abstraction. Answers questions like: how do I manage complexity, implement data structures, build abstractions, represent a domain, or encourage code reuse?

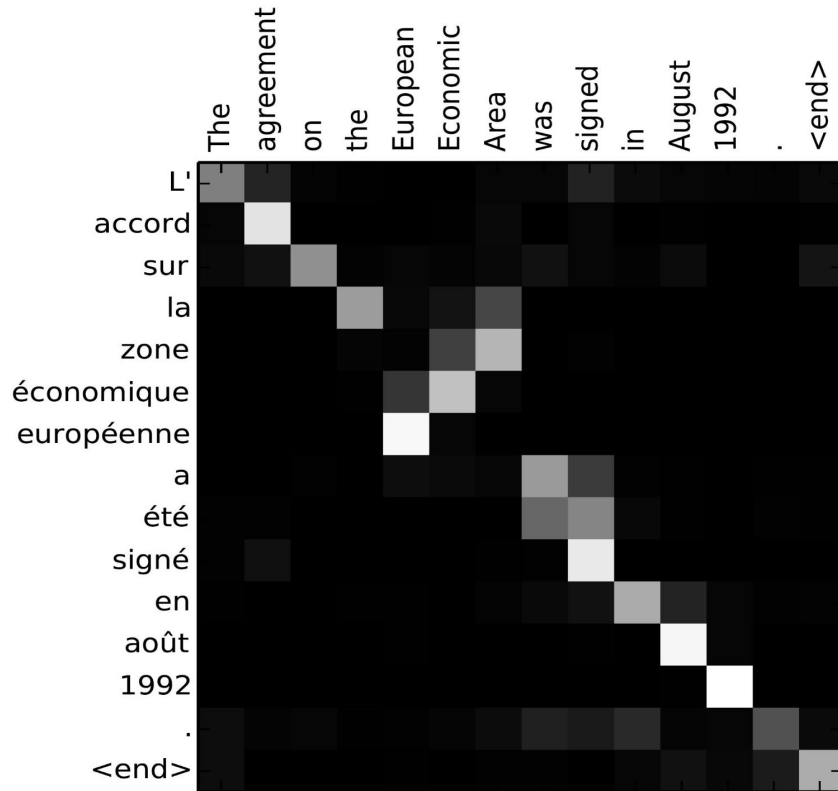
- Procedural programming – imperative procedures
- Object-oriented programming – objects
- Functional programming – pure functions
- Logic programming – predicates
- Concurrent programming – communicating processes
- Symbolic programming – rewrite rules
- Stack-oriented programming – stack operations
- Differentiable programming – ???

Layers

A lambda that closes over numerical parameters.

```
195
196 dense(W, b, σ = identity) =
197     x → σ.(W * x .+ b)
198
199 chain(f ... ) = foldl(∘, reverse(f))
200
201 mlp = chain(
202     dense(randn(5, 10), randn(5), tanh),
203     dense(randn(2, 5), randn(2)))
204
205 x = rand(10)
206
207 mlp(x)
208     ▾ Float64[2]
        0.646...
```

Data Structures: Attention



See also: Neural Turing Machine

Resources

- Julia Language: julialang.org
- Flux ML library: fluxml.ai
- Zygote AD: github.com/FluxML/Zygote.jl
- Differentiable Control fluxml.ai/blog/2019/03/05/dp-vs-rl.html
- Google Scholar for more papers and technical introductions