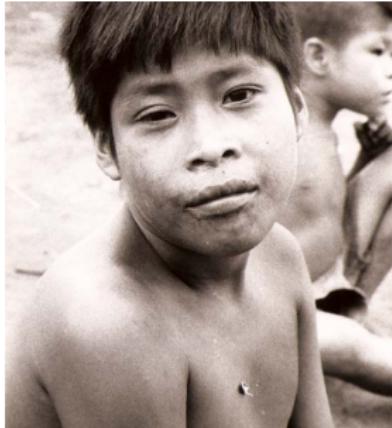


Reverse Mode Automatic Differentiation: Unraveling Expression Graphs & Library Magic

Steve Broder

October 2023



1978



2009



A



B

From:

Can Anthropologists Distinguish Good and Poor Hunters? Implications for Hunting Hypotheses, Sharing Conventions, and Cultural Transmission

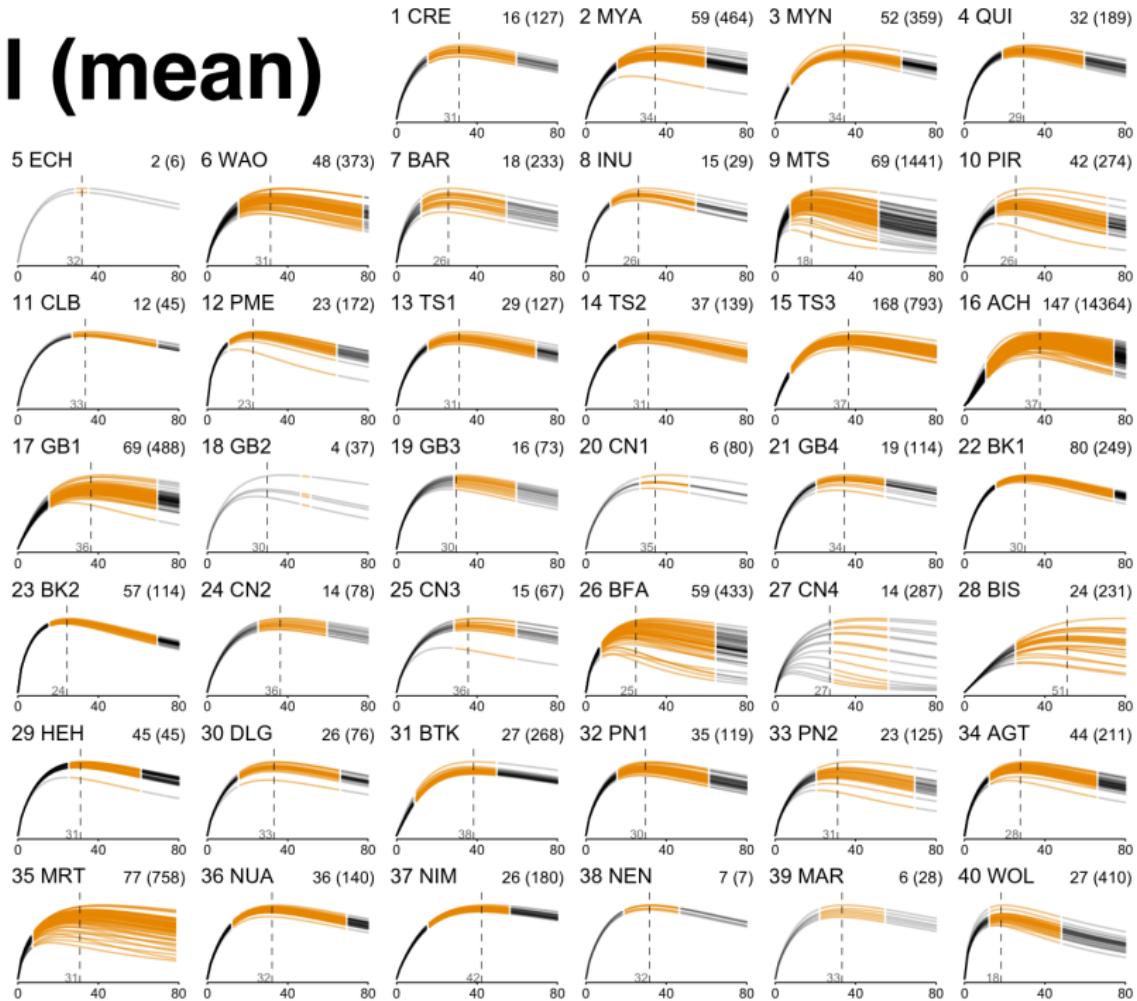
Kim Hill and Keith Kintigh

Life History of Production Skill

- Goals
 - Estimate skill development
 - Take variation seriously
 - Develop stats machinery
- Sample
 - 40 research sites
 - 1821 individuals
 - 21,160 foraging trips
 - 23,747 harvests
 - Uncountable headaches



Skill (mean)



Estimating COVID Infection Rates For Policy

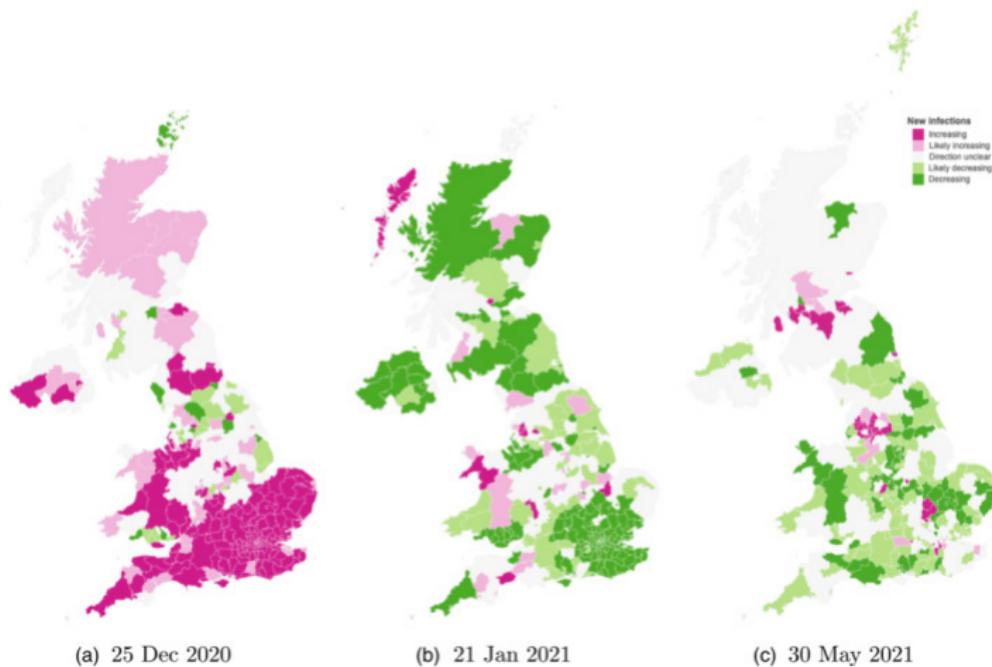


Figure: Probability of epidemic growth by local area

Automatic Differentiation Probably Effects Your Day to Day

ChatGPT 5 ▾

What programming technique is used to compute the gradients of your model during training?

The main programming technique used to compute gradients in models like mine (and in most modern deep learning frameworks) is **reverse-mode automatic differentiation (AD)**, also known as **backpropagation** when applied to neural networks.

What's Automatic Differentiation?

Computational technique for evaluating derivatives of functions expressed as computer programs by systematically applying the chain rule.



Figure: Asking ChatGPT to make a physical representation of automatic differentiation.

Why use Automatic Differentiation?

Ex: Newton's method for finding root of function (where output is 0)

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Why use Automatic Differentiation?

$$f(x) = x^3 + x^2 + x \quad (1)$$

$$f'(x) = 3x^2 + 2x + 1 \quad (2)$$

Why use Automatic Differentiation?

Think about HMC, BFGS, SGD, etc.

- ▶ HMC: $\frac{dp}{dt} = -\nabla_{\theta} \log p(\theta|y)$
- ▶ BFGS: $s_k = -H_k \nabla_{\theta} f(\theta_k)$
- ▶ SGD: $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t; x_t)$
- ▶ Choices
 - ▶ Write by hand
 - ▶ finite difference,
 - ▶ symbolic differentiation
 - ▶ spectral differentiation
 - ▶ automatic differentiation

Why use Automatic Differentiation?

$$\underbrace{p(\theta | y)}_{\text{posterior}} \propto \prod_{i=1}^N \left\{ \sum_{z_i \in \{1, 2, 3\}} T_i \underbrace{\left[\pi_{z_i, 1} \prod_{t=2}^{T_i} \Pi_{z_i, t-1, z_i, t} \right]}_{\text{3-state HMM prior}} \underbrace{\prod_{t=1}^{T_i} \mathcal{N}(y_{i,t} | \eta_{i,t}, \sigma_{z_i, t}^2)}_{\text{state-dependent emission}} \right\}$$

where $\eta_{i,t} = \underbrace{x_{i,t}^\top \beta}_{\text{fixed}} + \underbrace{z_{i,t}^{(G)\top} b_g[i] + z_{i,t}^{(C)\top} c_c[i] + u_g[i]}_{\text{crossed random effects}} + \underbrace{f(t_{i,t})}_{\text{GP}} + \underbrace{\mu_{z_i, t} + r_{z_i, t}^\top w_i}_{\text{state-specific offset + slope}}, \quad z_{i,t} \in \{1, 2, 3\}.$

$$p(f | \psi) = (2\pi)^{-T/2} |\mathbf{K}|^{-1/2} \exp\left(-\frac{1}{2} f^\top \mathbf{K}^{-1} f\right),$$

$$\mathbf{K} = \sigma_f^2 \left(\mathbf{K}_{LP}(\ell, p, \lambda) + \rho \mathbf{K}_{SE}(\tilde{\ell}) \right) + \sigma_n^2 \mathbf{I}, \quad [\mathbf{K}_{LP}]_{tt'} = \exp\left(-\frac{(t-t')^2}{2\ell^2} - \frac{2\sin^2(\pi|t-t'|/p)}{\lambda^2}\right),$$

$$[\mathbf{K}_{SE}]_{tt'} = \exp\left(-\frac{(t-t')^2}{2\tilde{\ell}^2}\right), \quad \mathbf{K} = \mathbf{L}_K \mathbf{L}_K^\top \Rightarrow \log |\mathbf{K}| = 2 \sum_{j=1}^T \log ((\mathbf{L}_K)_{jj}).$$

Hierarchical mixed effects (non-centered, LKJ prior):

$$\mathbf{b}_g = (\mathbf{I}_{p_G} \otimes \text{diag}(\boldsymbol{\tau}_b) \mathbf{L}_R) \tilde{\mathbf{b}}_g, \quad \tilde{\mathbf{b}}_g \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \mathbf{c}_c = \text{diag}(\boldsymbol{\tau}_c) \tilde{\mathbf{c}}_c, \quad \tilde{\mathbf{c}}_c \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

$$\text{LKJ}_{p_G}(\eta) \text{ prior on } \mathbf{R}, \quad \mathbf{L}_R \mathbf{L}_R^\top = \mathbf{R}, \quad \boldsymbol{\tau}_b \sim \prod_{j=1}^{p_G} \text{Half-}t_{\nu_b}(0, s_b), \quad \boldsymbol{\tau}_c \sim \prod_{j=1}^{p_C} \text{Half-}t_{\nu_c}(0, s_c),$$

$$u_g \sim \mathcal{N}(0, \sigma_u^2).$$

Why use Automatic Differentiation?

- ▶ Faster than finite difference, more flexible than symbolic differentiation
- ▶ Allows for unknown length while and for loops
- ▶ Accurate to floating point precision
- ▶ Reverse Mode AD can compute partials derivatives of inputs at the same time

How Fast is AutoDiff?



Figure: AuToDiFf rUnS iN $\Theta(C(f))$ TiMe

Impl Matters!

TODO: QR code here

Regression

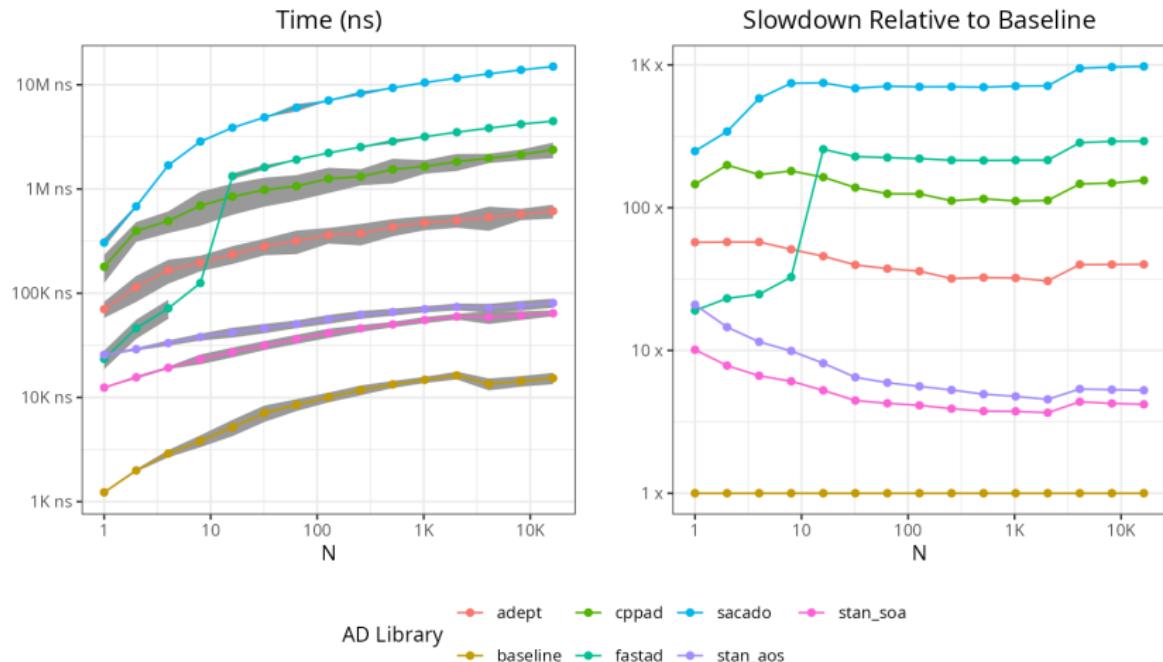


Figure: Benchmark for f given $f = y \sim N(X\theta, \sigma)$

Goal of this talk

- ▶ Understand performance of high throughput memory intensive programs
- ▶ Show how modern C++ through time has led to cleaner and more efficient AD

What's Automatic Differentiation?

- ▶ Given a function f with inputs $x \in \mathbb{R}^n$ and outputs $z \in \mathbb{R}^m$ we want to calculate the Jacobian J with size (m, n)
- ▶ To get the full Jacobian, use the chain rule to differentiate from each output to each input.

$$J_{i,1:j} = \left\{ \frac{\partial z_i}{\partial x_1}, \dots, \frac{\partial z_i}{\partial x_j} \right\}$$

Automatic Differentiation can do higher order partials, but here we just focus on the Jacobian

Cool Math, but how do we do this in a computer??

For Reverse Mode AD, we perform two functions.

- ▶ Forward Pass:

$$z = f(x_0, x_1)$$

- ▶ Reverse Pass: Given z 's adjoint (gradient) \bar{z}

$$\text{chain}(z, x_0, x_1) = \left\{ \frac{\partial z}{\partial x_0} \bar{z}, \frac{\partial z}{\partial x_1} \bar{z} \right\}$$

Calculate the adjoint-jacobian update for x_0 and x_1 .

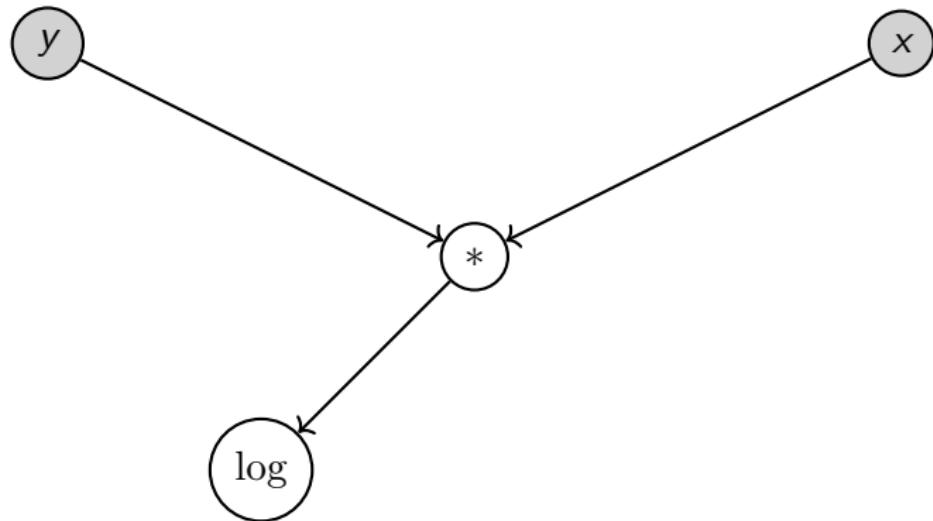
- ▶ The calculations needed are represented as an expression Graph

What's an Expression Graph?

- ▶ Dependency graph of intermediate computations
- ▶ Think of both data and operations as objects
- ▶ Do a forward pass to calculate the values of the intermediates, then a reverse pass to calculate the adjoint-jacobian updates.

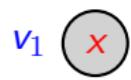
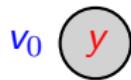
Forward Pass

$$z = \log(x * y)$$



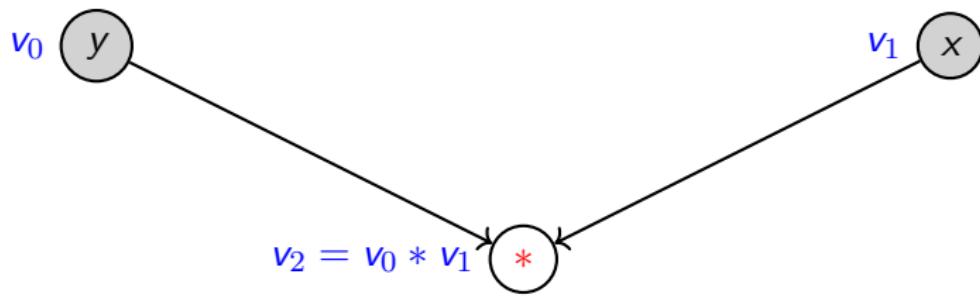
Forward Pass

$$z = \log(x * y)$$



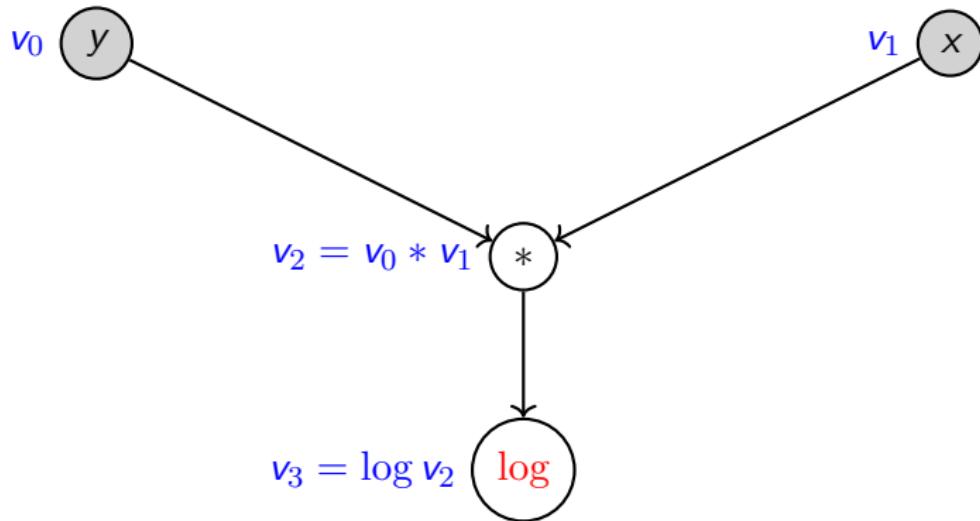
Forward Pass

$$z = \log(x * y)$$



Forward Pass

$$z = \log(x * y)$$



How do we calculate the adjoint jacobian?

Let \bar{v}_i be the adjoint of v_i

$$\bar{v}_i = \frac{\partial v_{i+1}}{\partial v_i} \bar{v}_{i+1}$$

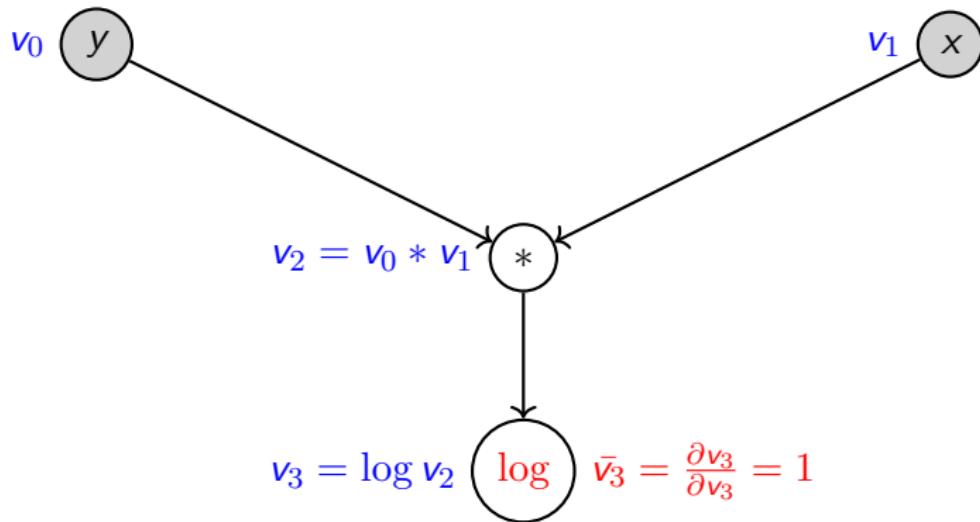
Automatic Differentiation only needs the partials of the intermediates

$$z = x * y \quad \frac{\partial z}{\partial x} = y, \frac{\partial z}{\partial y} = x$$

$$z = \log(x) \quad \frac{1}{x}$$

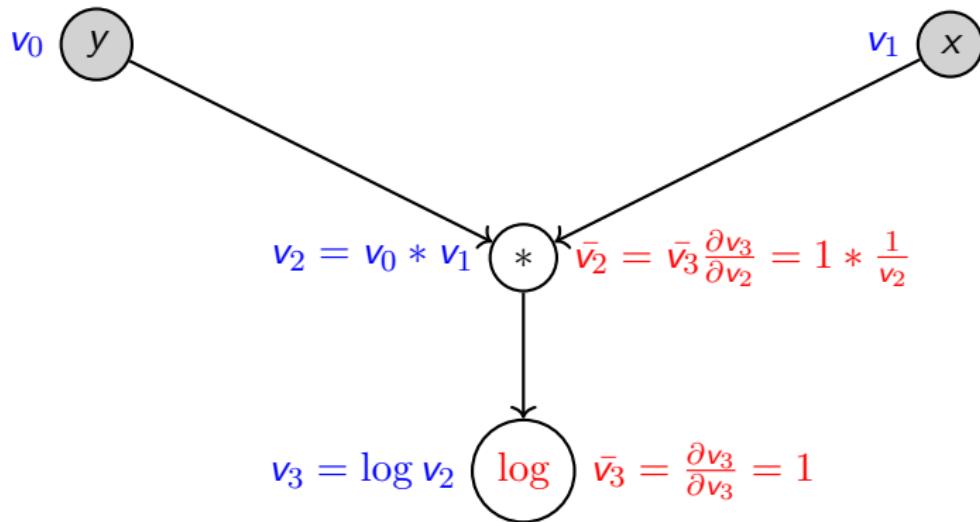
Reverse Pass

$$z = \log(x * y)$$



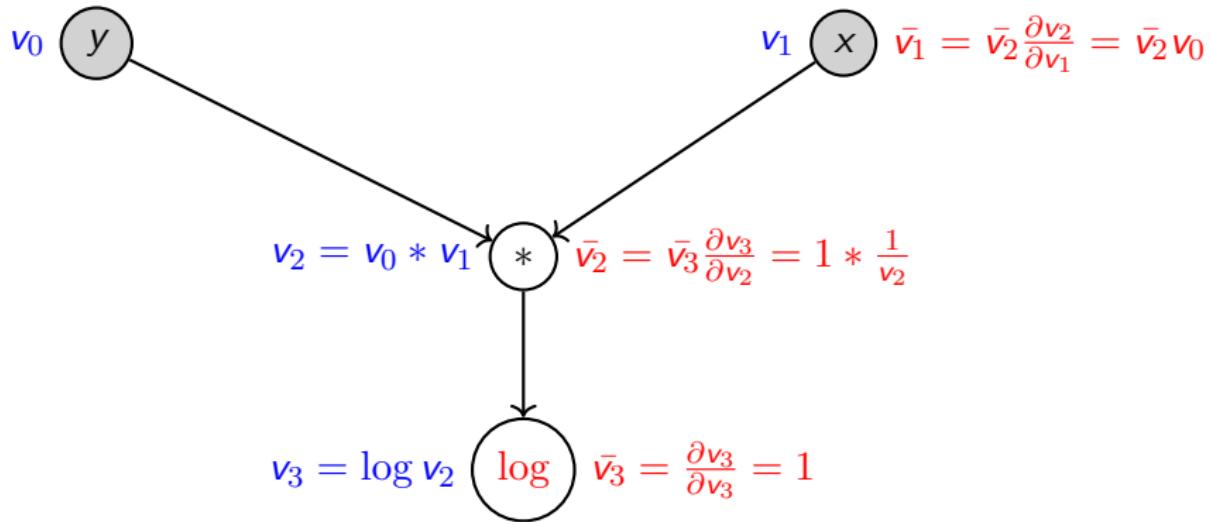
Reverse Pass

$$z = \log(x * y)$$



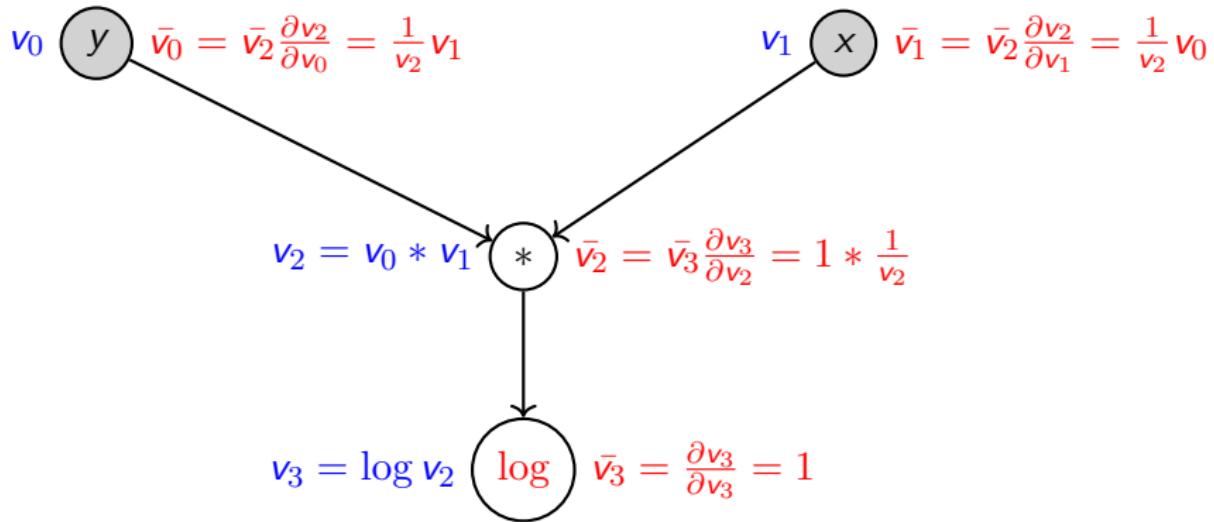
Reverse Pass

$$z = \log(x * y)$$



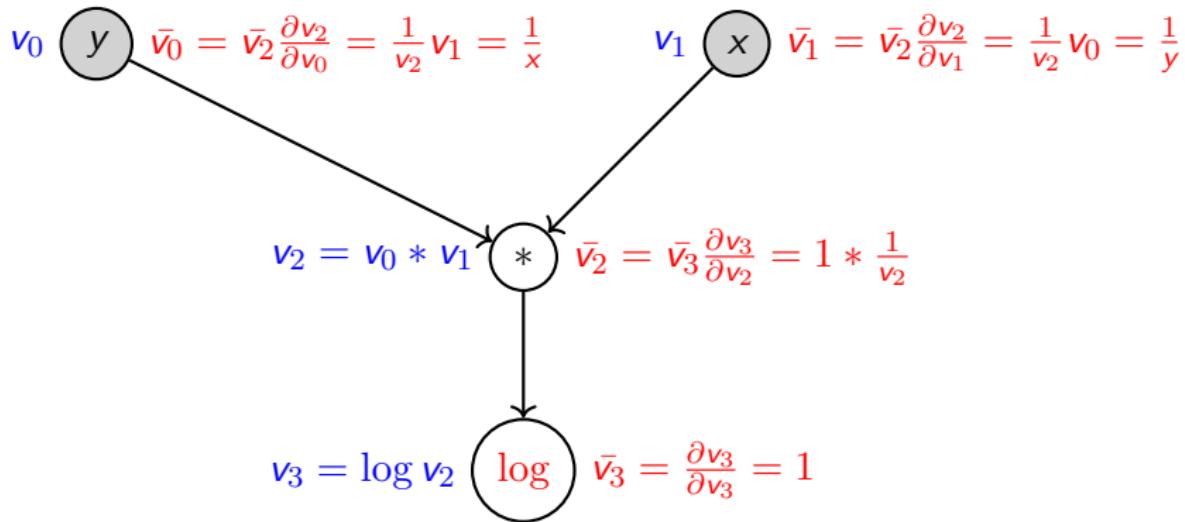
Reverse Pass

$$z = \log(x * y)$$



Reverse Pass

$$z = \log(x * y)$$



What Do AD Libraries Care About?

- ▶ Flexibility:
 - ▶ Debugging, exceptions, conditional loops, matrix subset assignment
- ▶ : Efficiency:
 - ▶ Efficiently using a single CPU/GPU
- ▶ Scaling
 - ▶ Efficiently using clusters with multi-gpu/cpu nodes

How do we keep track of our reverse pass?

- ▶ Source code transformation
 - ▶ Unroll all forward passes and reverse passes into one function
 - Good: Fast
 - Bad: Hard to implement, very restrictive

How do we keep track of our reverse pass?

- ▶ Source code transformation
 - ▶ Unroll all forward passes and reverse passes into one function
 - Good: Fast
 - Bad: Hard to implement, very restrictive
- ▶ Operator Overloading
 - ▶ Nodes in the expression graph are objects which store a forward and reverse pass function
 - Good: Easier to implement, more flexible
 - Bad: Less optimization opportunities

Newer AD packages use a combination of both

How do we keep track of our reverse pass?

Static (Fast) vs. Dynamic (Flexible) graphs

- ▶ Known expression graph size at compile time? (Static)
- ▶ Reassignment of variables (dynamic easy, Static vv hard!)
- ▶ How much time do I have? (dynamic)

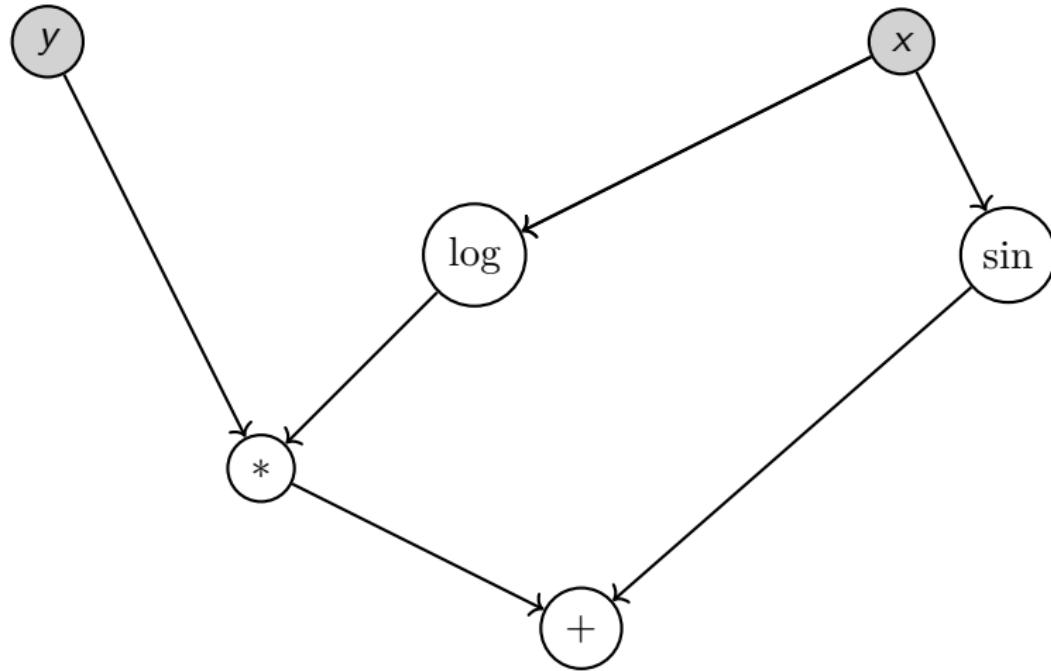
How do we keep track of our reverse pass?

Static (Fast) vs. Dynamic (Flexible) graphs

- ▶ Known expression graph size at compile time? (Static)
- ▶ Reassignment of variables (dynamic easy, Static vv hard!)
- ▶ How much time do I have? (dynamic)
- ▶ Do I actually need to track my graph :)

Make A Tape

$$z = \log(x) * y + \sin(x)$$



Cool graph math, how do we do this in a computer?

```
1 auto foo_fwd(double x, double y);
```

Cool graph math, how do we do this in a computer?

```
1 auto foo_fwd(double x, double y);  
2  
3 template<ADType Ret, ADType T1, ADType T2>  
4 auto foo_rev(Ret&& z, T1& x, T2& y) {  
5     adjoint(x) += adjoint(z) * compute_adj(x);  
6     adjoint(y) += adjoint(z) * compute_adj(y);  
7 }
```

Cool graph math, how do we do this in a computer?

```
1 auto foo_fwd(double x, double y);
2 template<ADType Ret, ADType T1, ADType T2>
3 auto foo_rev(Ret&& z, T1&& x, T2&& y) {
4     adjoint(x) += adjoint(z) * compute_adj(x);
5     adjoint(y) += adjoint(z) * compute_adj(y);
6 }
7 template<ADType T1, ADType T2>
8 auto foo(T1&& x, T2&& y) {
9     auto fwd_ret = foo_fwd(x, y);
10    auto foo_rev = [](auto&& z, auto& x, auto& y) {
11        return my_func_rev(z, x, y);
12    };
13    auto rev_ops = tuple(fwd_ret, forward_as_tuple(x,
14        ~ y));
15    return tuple{rev_ops, foo_rev};
16 }
```

Make A Tape

$$f(x, y) = \log(x)y + \sin(x)$$

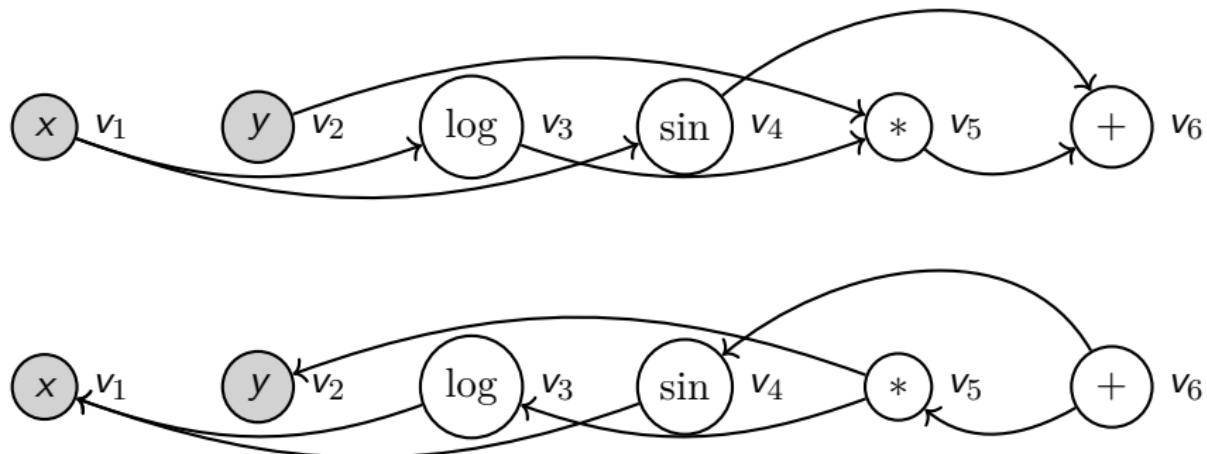


Figure: Topological sort of expression graph

Operator Overloading Approach

- ▶ The operator overloading approach usually involves:
 - ▶ Tracking the expression graph for the reverse pass function calls
 - ▶ A pair scalar type to hold the value and adjoint
- ▶ Very flexible: Allows conditional loops and reassignment of values in matrices
- ▶ Nodes of expression graph can be collapsed

[Example Godbolt](#)

Operator Overloading Approach: Adolc

```
1 struct vari {
2     double val;
3     double adj;
4 }; // 12 bytes
5 struct var {
6     vari* vi_;
7 }; // 8 bytes
8 struct ad_tape {
9     struct expr_info {
0         vari** out_;
1         vari** in_;
2         op_code chain_op;
3     };
4     std::vector<expr_info> tape_;
5     inline void add_expr(var out, var in1, var in2,
6         op_code expr_op) {
7         tape_.push_back({out.vi_, {in1.vi_, in2.vi_},
8             expr_op});
9         return out;
10    }
11};
```



Operator Overloading Approach: Adolc

```
1 inline auto operator+(var x, var y) {  
2     return tape.add_expr(var(x.val() + y.val()), x,  
3     → y, op_code::ADD_EXPR);  
4 }
```

Operator Overloading Approach

```
1 void calc_grad(Tape& tape) {
2     tape_val = tape.start();
3     while (tape.start() != tape.end()) {
4         auto op = tape.get_op();
5         switch (op.chain_op) {
6             case log_expr:
7                 op.in(0).adj += op.out(0).adj * 1.0 /
8                     op.in(0).val;
9             case mul_expr:
10                 op.in(0).adj += op.out(0).adj * op.in(1).val;
11                 op.in(1).adj += op.out(0).adj * op.in(0).val;
12         }
13     }
14 }
```

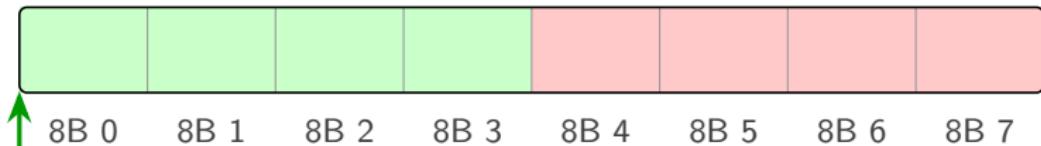
Operator Overloading Approach: Stan

```
1 struct vari {
2     double val;
3     double adj;
4     virtual void chain() {};
5 }; // 24 bytes
6 struct var {
7     vari* vi_;
8     var(double val) : vi_(new vari{val, 0.0}) {}
9 }; // 8 bytes
```



Operator Overloading Approach: Stan

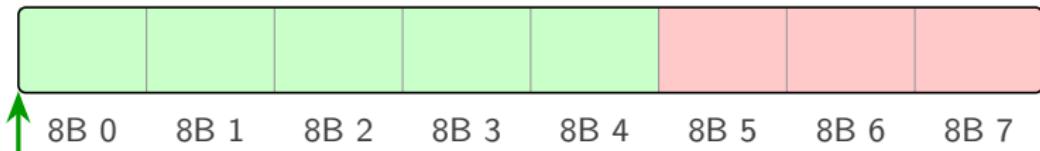
```
1 struct LogVari : vari {  
2     vari* in_;  
3     LogVari(var x) : in_(std::log(x.vi_->val)) {}  
4     virtual void chain() override {  
5         in_->adj += adj / in_->val;  
6     }  
7 }; // 32 bytes  
8 inline var log(var x) {  
9     return var(new LogVari(x));  
0 }
```



`LogVari{dbl, dbl, vptr, vari*}`

Operator Overloading Approach: Stan

```
1 struct add_vv : public var_impl {  
2     vari* lhs_;  
3     vari* rhs_;  
4     void chain() {  
5         lhs_->adjoint_ += this->adjoint_;  
6         rhs_->adjoint_+= this->adjoint_;  
7     }  
8 }; // 40 bytes
```



AddVari{dbl, dbl, vptr, vari*, vari*}

Operator Overloading Approach: Stan

```
1 var x = 1;  
2 var y = 2;  
3 var z = log(x * y);  
4 grad(z);
```

Operator Overloading Approach: Stan

```
1 template <typename T, typename F>
2 struct callback_vari : public vari_value<T> {
3     F rev_functor_;
4     template <typename S>
5     explicit callback_vari(S&& value, F&& rev_functor)
6         : vari_value<T>(std::move(value), true),
7             rev_functor_(std::forward<F>(rev_functor))
8             {}
9     inline void chain() final { rev_functor_(*this); }
0 }; // 24 + 8B * N
```

Operator Overloading Approach: Matrices

- ▶ Either Array of Structs (AOS) or Struct of Arrays (SoA)

```
1 struct var {
2     double value_;
3     double adjoint_;
4 };
5 struct MatrixVar {
6     var* data_;
7     MatrixVar(std::size_t N) :
8         data_(static_cast<var*>(malloc(sizeof(var) * N)))
9             {}}
10 MatrixVar aos_matrix;
11
12 struct VarMatrix {
13     double* value_;
14     double* adjoint_;
15 }
16 VarMatrix soa_matrix;
```

Source Code Transform Ex:

```
1 double z = log(x * y);
```

Break it down

```
1 double v0 = x;
2 double v1 = y;
3 double v2 = x * y;
4 double v3 = log(v2)
5 double bar_v3 = 1;
6 double bar_v2 = bar_v3 * 1/v2;
7 double bar_v1 = bar_v2 * v0;
8 double bar_v0 = bar_v2 * v1;
```

Source Code Transform Ex:

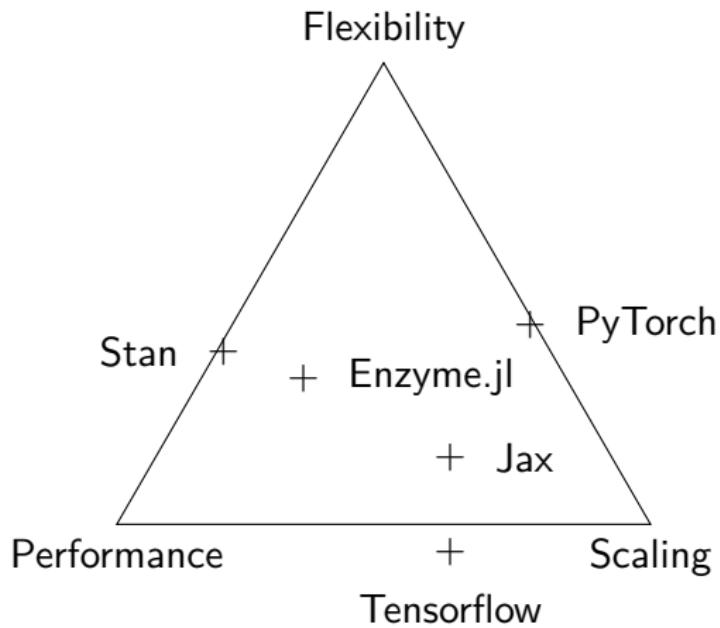
Code like the following very hard / impossible in source code transform

```
1 while(error < tolerance) {  
2     // ...  
3 }
```

Operator Overloading Approach: Matrices

- ▶ Array of Structs:
 - ▶ Simple, most algorithms Just Work™
 - ▶ Adds a lot to expression graph
 - ▶ turns off SIMD
- ▶ Struct of Arrays:
 - ▶ Hard, everything written out manually
 - ▶ Collapses matrix expressions in tree
 - ▶ SIMD can be used on values and adjoints

What are the AD packages like?



What are the AD packages like?

Disclaimer: Just pick the package that does the things you like,
the ones here are performant enough

Common Autodiff Packages

- ▶ Static Graph
 - ▶ TensorFlow, Jax, Enzyme
- ▶ Dynamic Graph
 - ▶ Pytorch and Stan
- ▶ TF, Jax, and Pytorch now have both

Stan!

Good:

- ▶ Very flexible language
- ▶ Exceptions, conditionals loops, matrix subsetting
- ▶ Only known CPU AD package faster than Stan math is [FastAD](#)
- ▶ Simple C like Domain Specific Language (DSL)

Bad:

- ▶ Very limited GPU support at the language level
- ▶ Poor scaling for TB of data
- ▶ Simple C like Domain Specific Language (DSL)
- ▶ Compilation times

Pytorch

Good:

- ▶ Good multi-gpu support
- ▶ Exceptions, conditional loops, debugging first priority
- ▶ Builtins for neural networks
- ▶ Extensible (see pytorch-finufft)

Bad:

- ▶ Subset assignment to matrices and vectors is a full hard copy
- ▶ Backend is very hard to parse

Tensorflow

Good:

- ▶ Made for scalability

Bad:

- ▶ No conditional loops
- ▶ No subset assignment to matrices and vectors
- ▶ No exceptions

Jax

Good:

- ▶ Built on top of autograd and XLA
- ▶ Well documented
- ▶ Extendable
- ▶ Write python, jit to near C++ speed

Bad:

- ▶ No Exceptions, conditional loops, subset assignment to matrices and vectors is a full hard copy

Enzyme.jl

Good:

- ▶ JIT compiled to llvm
- ▶ Can use a large amount of julia packages

Bad:

- ▶ Only one main maintainer
- ▶ Not yet 1.0 (0.1)
- ▶ No GC or dynamic dispatch support

What did we talk about?

- What's Automatic Differentiation (AD)?
 - ▶ Evaluates partial derivatives of a program
- Why should you care?
 - ▶ It's important and used a lot
- What's an expression graph?
 - ▶ Graphic to describe dependencies for AD
- How is AD implemented
 - ▶ Source code is transformed (static graph) or objects are made for intermediate ops (dynamic graph)
- What are the tradeoffs between different AD packages
 - ▶ Flexibility, Efficiency, and Scale