

+ 25

From Bayesian Inference to LLMs

Modern C++ Optimizations for Reverse-Mode Automatic Differentiation

STEVE BRONDER



20
25 | 
September 13 - 19

Goal of This talk

- Why you should care about Automatic Differentiation
- How Reverse Mode Automatic Differentiation works
- Performance of high throughput memory intensive programs
- How modern C++ has led to cleaner and more efficient AD

Estimating COVID Infection Rates For Policy

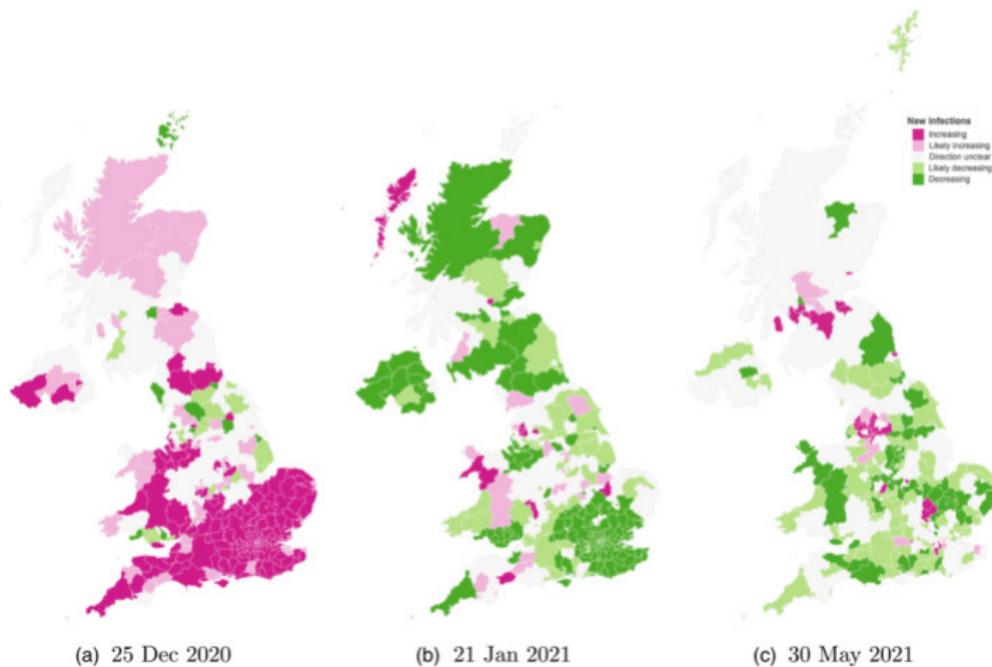


Figure: Probability of epidemic growth by local area

Show of Hands

Who has used an LLM?

What is Automatic Differentiation?

Automatically compute exact gradients of a function by treating a program as a graph of elementary operations and applying the chain rule on that graph

What is a gradient?

Ex: Newton's root finding method

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

find $f(x) = y$ where $y = 0$

What is a gradient?

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

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

Why use Automatic Differentiation?

- Hamiltonian Monte Carlo:

$$\frac{dp}{dt} = \nabla_{\theta} \log p(\theta|y)$$

- BFGS:

$$s_k = -H_k \nabla_{\theta} f(\theta_k)$$

- Stochastic Gradient Descent:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t; x_t)$$

Why use Automatic Differentiation?

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?

- I do not wish to write those gradients by hand

Why use Automatic Differentiation?

- I do not wish to write those gradients by hand
- Instead of N finite differences, a forward and backward pass

Why use Automatic Differentiation?

- I do not wish to write those gradients by hand
- Instead of N finite differences, a forward and backward pass
- Accurate to floating point precision

Why use Automatic Differentiation?

- I do not wish to write those gradients by hand
- Instead of N finite differences, a forward and backward pass
- Accurate to floating point precision
- Allows for unknown length while and for loops

Implementation Matters!

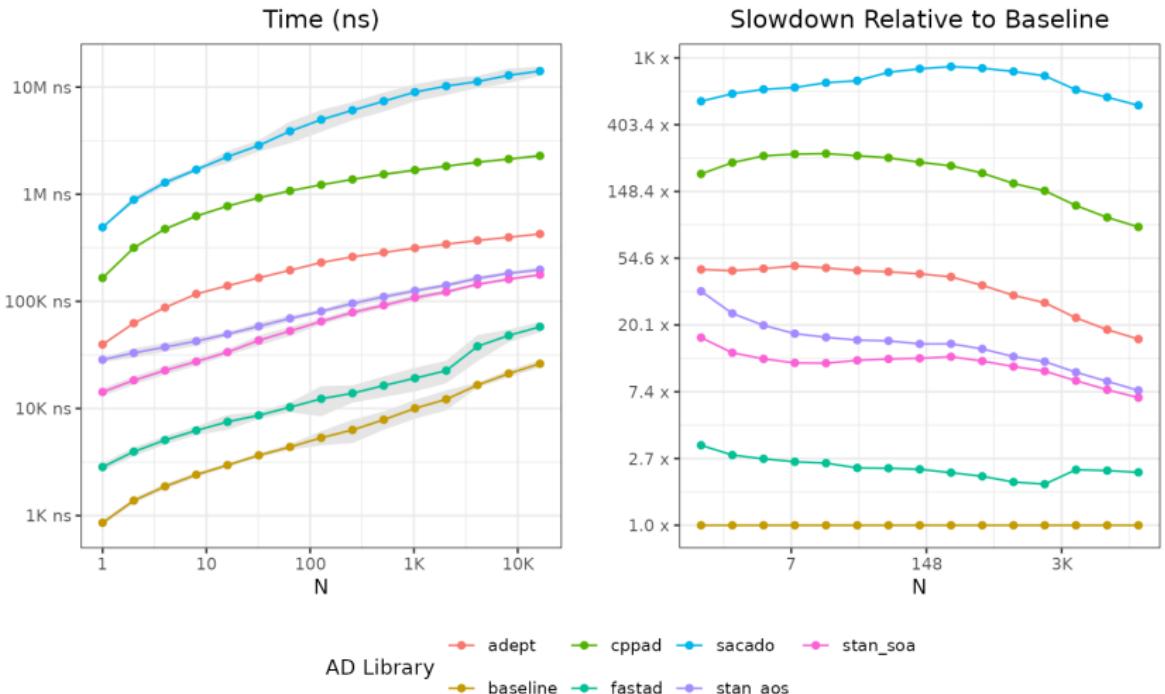


Figure: Log-Log Benchmark for f and f' given $f(y) = N(y|X\theta, \sigma)$

What is Automatic Differentiation?

- AD computes gradients of a program by applying the chain rule to its subexpressions.

$$f(x, y) = \log(x \cdot y)$$

$$g(x, y) = x \cdot y$$

$$h(u) = \log(u)$$

$$f(x, y) = h(g(x, y))$$

What is Automatic Differentiation?

- AD computes gradients of a program by applying the chain rule to its subexpressions.

$$f(x, y) = h(g(x, y))$$

$$g(x, y) = x \cdot y \quad g_x'(x, y) = y \quad g_y'(x, y) = x$$

$$h(u) = \log(u) \quad h'_u(u) = u^{-1}$$

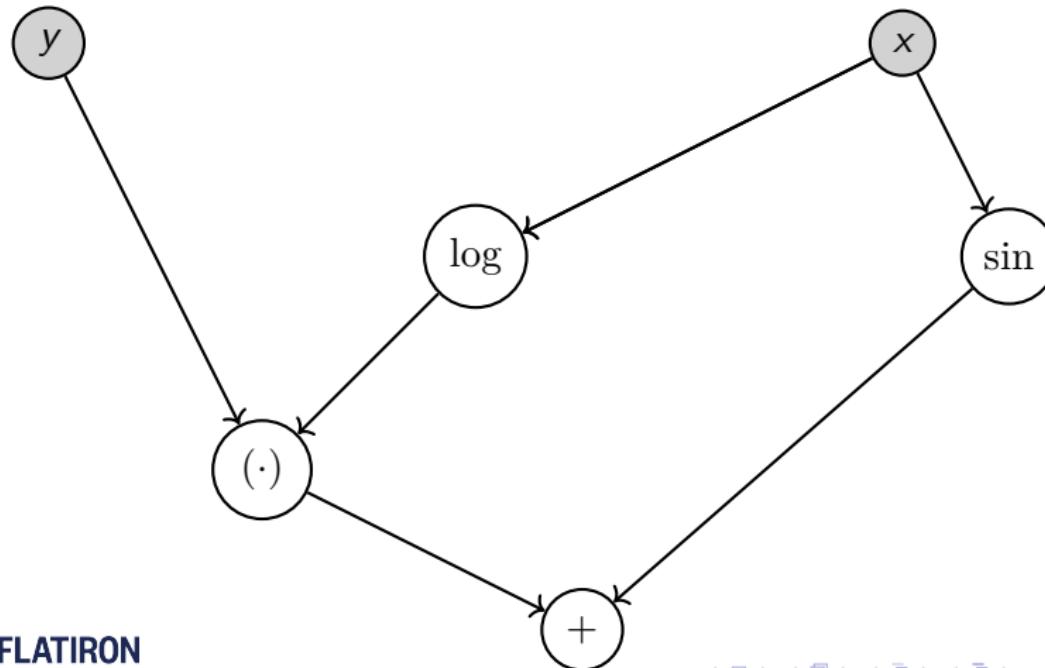
$$\begin{aligned} f'_x(x, y) &= h'_g(g(x, y)) \cdot g'_x(x, y) \\ &= \frac{1}{g(x, y)} \cdot y = \frac{1}{x}, \end{aligned}$$

$$\begin{aligned} f'_y(x, y) &= h'_g(g(x, y)) \cdot g'_y(x, y) \\ &= \frac{1}{g(x, y)} \cdot x = \frac{1}{y}. \end{aligned}$$

Data Type: Expression Graph

Goal: Calculate the full gradient by accumulating partial gradients (adjoints) through the a graph of subexpressions.

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



Expression Graph: Traversal

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

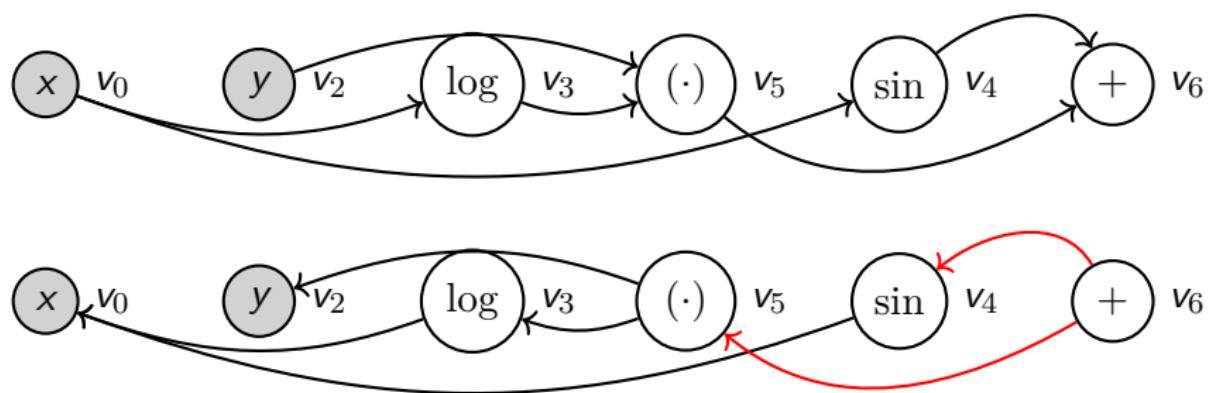


Figure: Topological sort of expression graph

Expression Graph: Node

For Reverse Mode AD, each node performs two functions.

- Forward Pass:

$$z_{child} = f(z_{parent_1}, z_{parent_2})$$

Expression Graph: Node

For Reverse Mode AD, each node performs two functions.

- Forward Pass:

$$z_{child} = f(z_{parent_1}, z_{parent_2})$$

- Reverse Pass:

Given z_{child} 's adjoint (partial gradient) \bar{z}_{child}

Calculate the local adjoint update for z_{parent_1} and z_{parent_2} .

$$chain(\bar{z}_{child}, z_{parent_1}, z_{parent_2}) = \begin{cases} \bar{z}_{parent_1} += \bar{z}_{child} \cdot z'_{parent_1} \\ \bar{z}_{parent_2} += \bar{z}_{child} \cdot z'_{parent_2} \end{cases}$$

Pseudocode for AutoDiff

```
// compute and store node.value
struct node {
    double val;
    double adj{0};
    virtual void chain() {}
};

node x = 2.0;
node y = 3.0;
// Forward pass calculated here
auto expr_graph = log(x) * y + sin(x);
// Calculate reverse pass
e_graph.end().adj = 1.0;
for (auto& expr : expr_graph | std::views::reverse) {
    expr->chain();
}
assert(y.adj == 1.2103);
assert(x.adj == 0.2516);
```

What Do AD Libraries Care About?

- Flexibility:
 - Debugging, exceptions, conditional loops, matrix 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?

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

Source code transformation

- Unroll all forward passes and reverse passes into one function

Good: Fast

Bad: Hard to implement, very restrictive

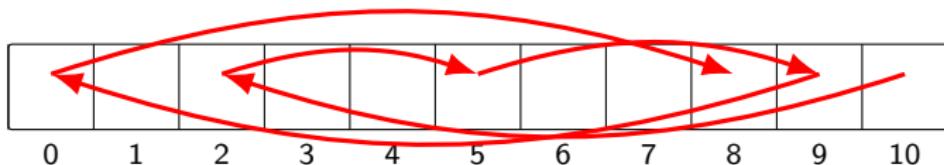
Operator Overloading Approach

The operator overloading approach usually involves:

- Pair scalar type to hold the value and adjoint
- Runtime "tape" for tracking expressions
- Arena allocator for storing AD nodes

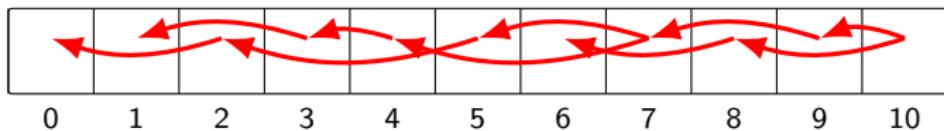
Issues

Noncontiguous Tape



Issues

Contiguous Tape



Operator Overloading: Monotonic Buffer

```
struct var_impl;
struct Tape {
    std::monotonic_buffer_resource mbr{(1 << 16)};
    std::polymorphic_allocator<std::byte> pa{&mbr};
    std::vector<var_impl*> tape;
    void clear() {
        tape.clear();
        mbr.release();
    }
};

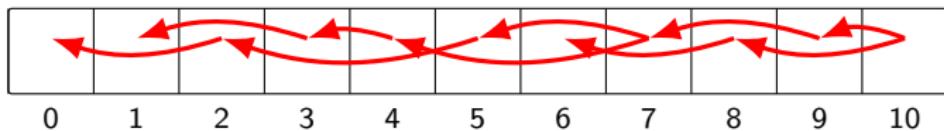
static Tape g_tape{};
```

Operator Overloading: Basic Type

```
struct var_impl {
    double val;
    double adj;
    virtual void chain() {};
}; // 24 bytes
struct var {
    var_impl* vi_;
}; // 8 bytes
```

Issues

Contiguous Tape



Operator Overloading: Expression Node

```
struct mul_vv : public var_impl {
    var op1_;
    var op2_;
    void chain() {
        op1_.adj += this->adjoint_ * op2_.val;
        op2_.adj += this->adjoint_ * op1_.val;
    }
}; // 40 bytes
```

Operator Overloading: Monotonic Buffer

```
template <typename T, typename Types>
auto new_node(Types&&... args) {
    auto* ret = g_tape.pa.template
        ↵ new_object<T>(args...);
    g_tape.tape.push_back(ret);
    return var{ret};
}
```

Operator Overloading: Expression Node

```
var operator*(var op1, var op2) {  
    return new_node<mul_vv>(op1.val() * op2.val(),  
        op1, op2);  
}
```

Operator Overloading: Reverse Pass

```
void grad(var& v) {
    v.adj = 1.0;
    for (auto&& expr : tape | std::views::reverse) {
        expr->chain();
    }
}
```

Operator Overloading: Example

```
void compute_grads(var x, var y) {
    var z = -10;
    while (z.val() < 20) {
        z += log(x * y);
    }
    grad(z);
}
```

Operator Overloading: Example

```
// Later in program
for (int i = 0; i < 1e10; ++i) {
    var x = compute_x(...);
    var y = compute_y(...);
    compute_grads(x, y);
    do_something_with_grads(x.adj, y.adj);
    g_tape.clear();
}
```

So Many Operator Classes

```
struct mul_vv;
struct mul_vd;
struct mul_dv;
struct add_vv;
struct add_dv;
struct add_vd;
struct subtract_vv;
struct subtract_dv;
struct subtract_vd;
struct divide_vv;
struct divide_dv;
struct divide_vd;
```

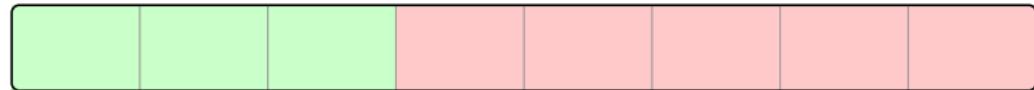
Reduce Boilerplate

```
template <Functor F>
struct callback_var_impl : public var_impl {
    F rev_functor_;
    template <std::floating_point S>
    explicit callback_var_impl(
        S&& value, F&& rev_functor)
        : var_impl(value),
        rev_functor_(rev_functor) {}
    void chain() final { rev_functor_(*this); }
}; // 24 + 8B * N
```

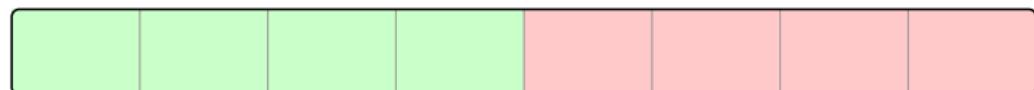
Reduce Boilerplate

```
template <typename T1, typename T2>
requires any_var<T1, T2>
inline auto operator*(T1 op1, T2 op2) {
    return new_node<callback_var_impl>(
        value(op1) * value(op2),
        [op1, op2](auto&& ret) mutable {
            if constexpr (is_var_v<T1>) {
                adjoint(op1) += adjoint(ret) * value(op2);
            }
            if constexpr (is_var_v<T2>) {
                adjoint(op2) += adjoint(ret) * value(op1);
            }
        });
}
```

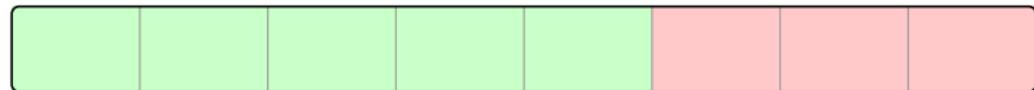
Poor Cache Use



`var_impl{dbl, dbl, vptr}`



`unary_var{dbl, dbl, vptr, var_impl*}`



`binary_var{dbl, dbl, vptr, var_impl*, var_impl*}`

Source Code Transform Example:

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

Break it down

```
double v0 = x;
double v1 = y;
double v2 = x * y;
double v3 = log(v2)
double bar_v3 = 1;
double bar_v2 = bar_v3 * 1/v2;
double bar_v1 = bar_v2 * v0;
double bar_v0 = bar_v2 * v1;
```

Source Code Transform Base Type:

```
struct var {
    double values_;
    double adjoints_;
    var(double x) : values_(x), adjoints_(0) {}
};
```

Source Code Transform Expression Node

```
template <typename F, typename... Exprs>
struct expr {
    var ret_;
    std::tuple<deduce_ownership_t<Exprs>...> exprs_;
    std::decay_t<F> f_;
    template <typename FF, typename... Args>
    expr(double x, FF&& f, Args&&... args) :
        ret_(x), f_(std::forward<F>(f)),
        exprs_(std::forward<Args>(args)...)
};
```

Source Code Transform Operator

```
template <typename T1, typename T2>
requires any_var_or_expr<T1, T2>
inline auto operator*(T1&& lhs, T2&& rhs) {
    return make_expr(value(lhs) * value(rhs),
        [](auto&& ret, auto&& lhs, auto&& rhs) {
            if constexpr (is_var_v<T1>) {
                adjoint(lhs) += adjoint(ret) * value(rhs);
            }
            if constexpr (is_var_v<T2>) {
                adjoint(rhs) += adjoint(ret) * value(lhs);
            }
        }, std::forward<T1>(lhs), std::forward<T2>(rhs));
}
```

Source Code Transform Expression Graph

```
auto z = x * log(y) + log(x * y) * y;  
// Polish Notation  
// + (* x (log y)) (* (log (* x y)) y)  
expr<Lambda<Plus>,  
    expr<Lambda<Mult>, var, expr<Lambda<Log>, var>>,  
    expr<Lambda<Mult>,  
    expr<Lambda<Log>, expr<Lambda<Mult>, var, var>>,  
    var>>
```

Source Code Transform Gradient

```
template <typename Expr>
inline void grad(Expr&& z) {
    adjoint(z) = 1.0;
    // Tuple nodes to be iterated in breadth-first
    // → order
    // <expr<Plus, expr<Mult>, expr<Mult>>,
    // expr<Mult, var, expr<Log>>,
    // expr<Mult, expr<Log>, expr<Mult>> ...
    auto nodes = collect_bfs(z);
    eval_breadthwise(nodes);
}
```

Example Godbolt

Scalar Comparison

Table: $f(x, y) = x \log(y) + \log(xy)y;$

Method	CPU Time	Slowdown
Dynamic	47.8ns	9x
Source Code Transform	6.3ns	.21x
Baseline	5.23ns	

Source Code Transform: Downsides

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

```
while(error < tolerance) {  
    var z += log(x * y);  
}
```

Operator Overloading Approach: Matrices

```
Matrix<var> B(M, M);
Matrix<var> X(M, M);
Matrix<var> Z = X * B.transpose();
```

Operator Overloading Approach: Matrices

```
template <typename MatrixType>
struct arena_matrix :  
    public Eigen::Map<MatrixType> {  
    using Base = Eigen::Map<MatrixType>;  
    template <typename T>  
    arena_matrix(T&& mat) :  
        Base(copy_to_arena(mat.data(), mat.size()),  
              mat.rows(), mat.cols()) {}  
};
```

Operator Overloading Approach: Matrices

```
// Array of Structs
struct Matrix<var> {
    var* data_;
};

// Struct of Arrays
struct var_impl<Matrix<double>> {
    arena_matrix<double> value_;
    arena_matrix<double> adjoint_;
    virtual void chain() {}
};
```

Operator Overloading Approach: Matrices

```
// Array of Structs
struct Matrix<var> {
    var* data_;
};
```

- Array of Structs:

- Simple, most algorithms Just Work™

- Adds a lot to expression graph

- turns off SIMD

Operator Overloading Approach: Matrices

```
// Struct of Arrays
struct var_impl<Matrix<double>> {
    arena_matrix<double> value_;
    arena_matrix<double> adjoint_;
    virtual void chain() {}
};
```

- Struct of Arrays:

- Hard, everything written out manually

- Collapses matrix expressions in tree

- SIMD can be used on values and adjoints

Matrix Multiplication Example

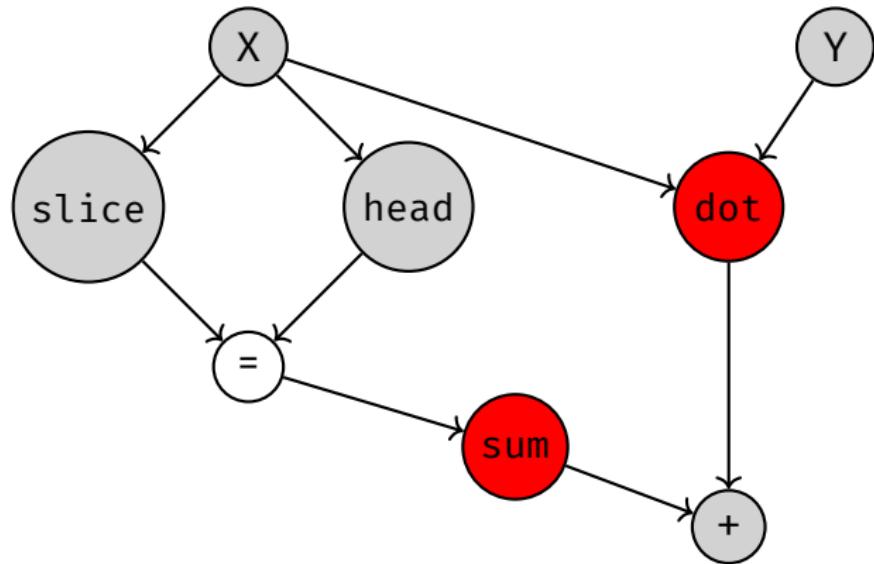
```
template <typename T1, typename T2>
requires any_var_matrix<T1, T2>
inline auto operator*(T1&& op1, T2&& op2) {
    return lambda_var(value(op1) * value(op2),
        [op1, op2](auto&& ret) mutable {
            if constexpr (is_var_matrix_v<T1>) {
                adjoint(op1) += adjoint(ret) *
                    → value(op2).transpose();
            }
            if constexpr (is_var_matrix_v<T2>) {
                adjoint(op2) += value(op1) * adjoint(ret);
            }
        });
}
```

Subset Assignment

```
var<Vector<double>> y{{0, 1, 2, 3}};
var<Vector<double>> x{{0, 1, 2, 3}};
var prod = y.dot(x);
x.slice(1, 3) = x.head(3);
auto z = prod + sum(x);
```

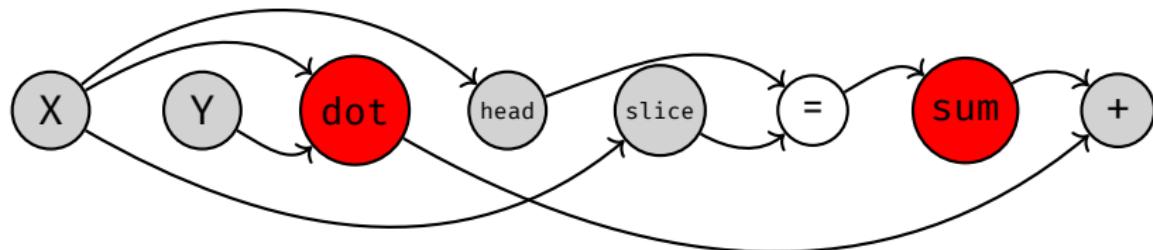
Subset Assignment

```
var<Vector<double>> y{{0, 1, 2, 3}};
var<Vector<double>> x{{0, 1, 2, 3}};
var prod = y.dot(x);
x.slice(1, 3) = x.head(3);
auto z = prod + sum(x);
```



Subset Assignment

```
var<Vector<double>> y{{0, 1, 2, 3}};
var<Vector<double>> x{{0, 1, 2, 3}};
var prod = y.dot(x);
x.slice(1, 3) = x.head(3);
auto z = prod + sum(x);
```



Subset Assignment Becomes Very Hard

```
var<Vector<double>> x{{0, 1, 2, 3}};  
x.slice(1, 3) = x.head(3);
```

Iter	X
init	{0, 1, 2, 3}
0	{0, 0, 2, 3}
1	{0, 0, 0, 3}
2	{0, 0, 0, 3}

Subset Assignment Becomes Very Hard

```
template <typename T>
struct var {
    template <typename S>
    require AssignableExpression<T, S>
    inline var<T>& operator=(const var<S>& other) {
        arena_matrix<T> prev_val(vi_->val_.rows(),
        ↳ vi_->val_.cols());
        prev_val.deep_copy(vi_->val_);
        vi_->val_.deep_copy(other.val());
        g_tape.callback(
            [this_vi = this->vi_, other_vi = other.vi_,
            ↳ prev_val]() mutable {
                this_vi->val_.deep_copy(prev_val);
                prev_val.deep_copy(this_vi->adj_);
                this_vi->adj_.setZero();
                other_vi->adj_ += prev_val;
            });
        return *this;
    }
}
```

Source Code Transform: Matrices

```
ad::Var X(4, 4); // Placeholders
ad::Var Y(4, 4);
auto expr = sum(X * Y);
auto [val_size, adj_size] = expr.cache_sizes();
Eigen::VectorXd val_buf(val_size);
Eigen::VectorXd adj_buf(adj_size);
expr.bind(std::make_pair(val_buf, 0),
         std::make_pair(adj_buf, 0));
```

Matrix Expression Graph Gradient

```
template <typename Expr>
void grad(Expr&& expr,
          std::pair<MatrixXd, MatrixXd>& X,
          std::pair<MatrixXd, MatrixXd>& Y) {
    expr.reset_adjoint();
    ad::autodiff(expr, ad::VarView(X.first, X.second),
                  ad::VarView(Y.first, Y.second));
}
```

Placeholder Node

```
template <typename T>
struct Var {
    double* val_;
    double* adj_;
    std::size_t rows_;
    std::size_t cols_;
    static constexpr std::size_t ops = 1;
    Var(std::size_t rows, std::size_t cols)
        : rows_(rows), cols_(cols) {}
    inline auto Forward(VarView<T>& x) {
        val_ = x.val_.data();
        adj_ = x.adj_.data();
        return Map<MatrixXd>(val_, rows_, cols_);
    }
    template <typename Arg>
    inline void Reverse(Arg& child_adj) {
        Map<MatrixXd>(adj_, rows_, cols_) += child_adj;
    }
}
```

Matrix Multiplication Expression Node

```
template <Matrix Op1, Matrix Op2>
struct MatMul {
    static constexpr std::size_t ops = 2;
    MatMul(const Op1& op1, const Op2& op2) :
        op1_(op1), op2_(op2) {}
    template <typename... Args>
    inline auto Forward(Args&&... args);
    template <typename Arg>
    inline auto Reverse(Arg&& child_adj);
    std::pair<std::size_t, std::size_t> cache_sizes();
    void bind_cache(double* val_buf, double* adj_buf,
                    std::size_t val_pos, std::size_t adj_pos);
    Op1 op1_;
    Op2 op2_;
    double* val_;
    double* adj_;
    std::size_t rows_;
    std::size_t cols_;
};
```

Expression Forward Pass

```
template <typename... Args>
inline auto Forward(Args&&... args) {
    auto op1_v = op1_.Forward(
        slice<0, Op1::ops>(args)...);
    auto op2_v = op2_.Forward(
        slice<Op1::ops, Op2::ops>(args)...);
    return this->val_ = op1_v * op2_v;
}
```

Expression Reverse Pass

```
template <Matrix Op1, Matrix Op2>
struct MatMul {
    template <typename Arg>
    inline auto Reverse(Arg&& child_adj) {
        this->adj += child_adj;
        auto left_adj = op1_.transpose() * this->adj;
        op2_.Reverse(left_adj);
        auto right_adj = this->adj * op2_.transpose();
        op1_.Reverse(right_adj);
    }
};
```

Expression Cache Size

```
template <Matrix Op1, Matrix Op2>
struct MatMul {
    inline auto cache_bind_size(
        std::pair<size_t, size_t>& child_sizes) {
        auto op1_sizes = op1_.cache_sizes();
        auto op2_sizes = op2_.cache_sizes();
        return std::make_pair(
            op1_sizes.first + op2_sizes.first +
            ↪ this->val_.size(),
            op1_sizes.second + op2_sizes.second +
            ↪ this->adj_.size());
    }
};
```

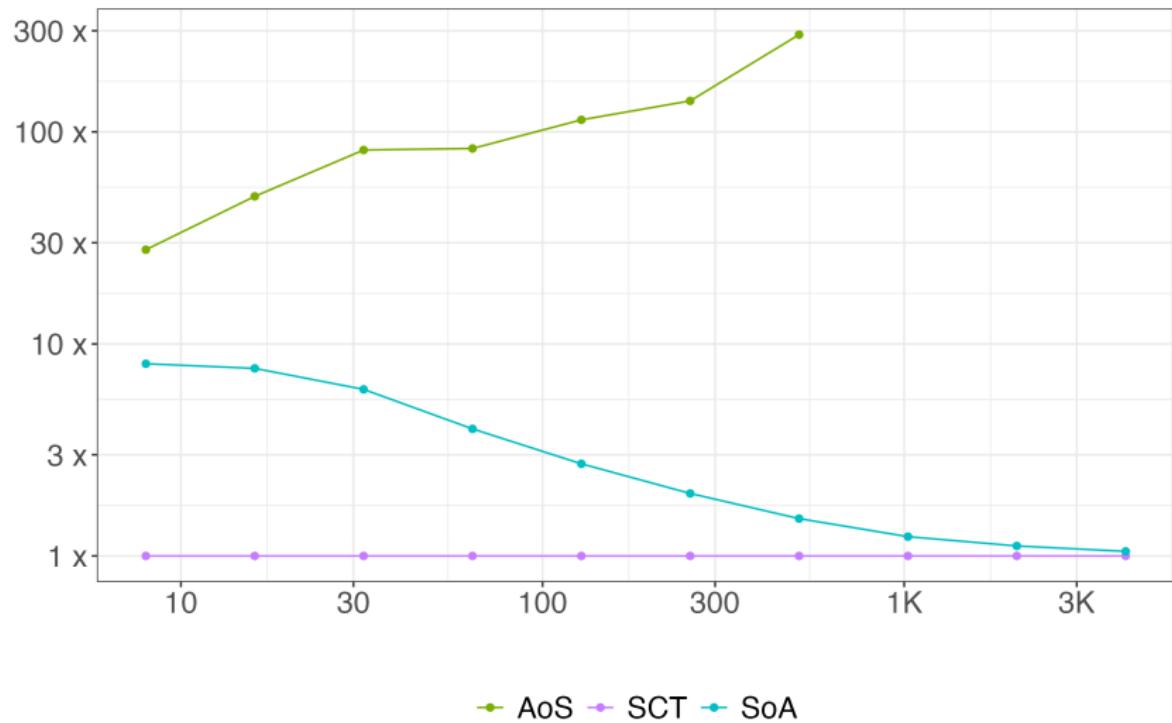
Expression Cache Size

```
template <Matrix Op1, Matrix Op2>
struct MatMul {
    void bind_cache(
        std::pair<double, size_t>& val_buf,
        std::pair<double*, size_t>& adj_buf) {
        this->val_ = val_buf.first + val_buf.second;
        this->adj_ = adj_buf.first + adj_pos.second;
        val_buf.second += this->val_.size();
        adj_buf.second += this->adj_.size();
        op1_.bind_cache(val_buf, adj_buf);
        op2_.bind_cache(val_buf, adj_buf);
    }
};
```

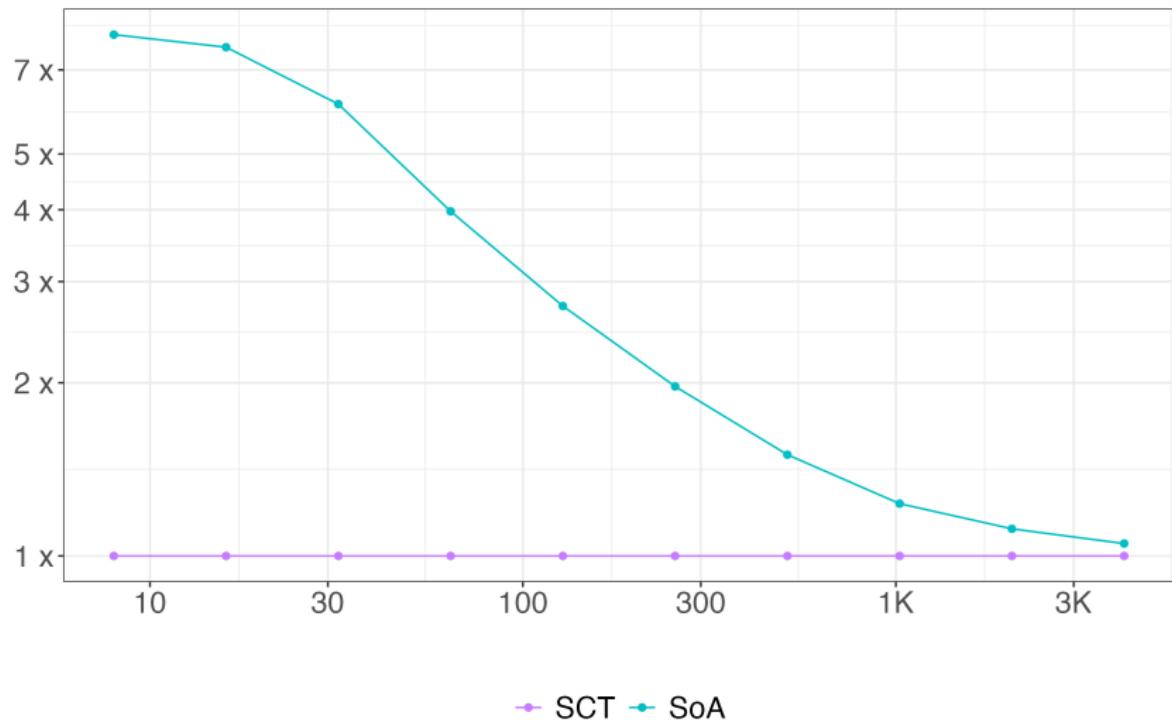
Matrix Expression Graph Creation

```
ad::Var X(4, 4); // Placeholders
ad::Var Y(4, 4);
auto expr = sum(X * Y);
auto [val_size, adj_size] = expr.cache_sizes();
Eigen::VectorXd val_buf(val_size);
Eigen::VectorXd adj_buf(adj_size);
expr.bind(std::make_pair(val_buf, 0),
          std::make_pair(adj_buf, 0));
ad::autodiff(expr,
              ad::VarView(X_v.first, X_v.second),
              ad::VarView(Y_v.first, Y_v.second));
```

Matrix Multiplication Benchmark



Matrix Multiplication Benchmark



— ● SCT — ● SoA

Speedup Holds for Larger Expression Graphs

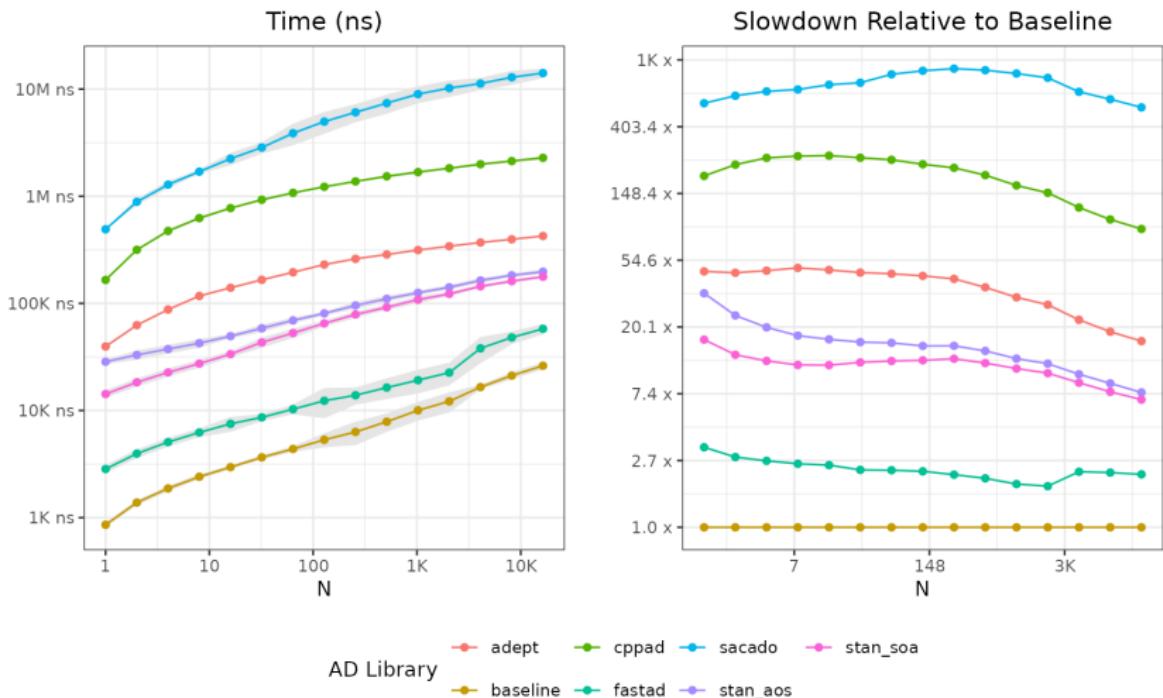


Figure: Log-Log Benchmark for f and f' given $f(y) = N(y|X\theta, \sigma)$

Thanks!

Repository for benchmarks and slides

