# Taking Derivatives of Functional Programs
## AD in a Functional Framework

Jeffrey Mark Siskind
qobi@purdue.edu

School of Electrical and Computer Engineering
Purdue University

University of Hertfordshire
22 November 2005

Joint work with Barak A. Pearlmutter.

## Outline

# Outline

1. **Lambda Calculus**

2. Differential Calculus in Lambda-Calculus Notation

3. Tutorial on AD
   - Forward Mode
   - Reverse Mode

4. Essence of the Derivation of Functional Reverse Mode

5. AD in Lambda-Calculus Notation

6. Examples

7. Benefits of this Approach

## Higher-Order Functions

$$\text{FOLD } (m, n, u, b, i) \stackrel{\triangle}{=} \textbf{if } m > n$$
$$\textbf{then } i$$
$$\textbf{else } b\,((u\ m), (\text{FOLD } (m + 1, n, u, b, i)))$$

# Higher-Order Functions

$$\text{FOLD } (m, n, u, b, i) \stackrel{\triangle}{=} \textbf{if } m > n$$
$$\textbf{then } i$$
$$\textbf{else } b \left((u\ m), (\text{FOLD } (m + 1, n, u, b, i))\right)$$

$$\sum_{i=m}^{n} \sin\ i$$

## Higher-Order Functions

$$\text{FOLD } (m, n, u, b, i) \stackrel{\triangle}{=} \textbf{if } m > n$$
$$\textbf{then } i$$
$$\textbf{else } b \left( (u\ m), (\text{FOLD } (m + 1, n, u, b, i)) \right)$$

$$\sum_{i=m}^{n} \sin\ i : \text{FOLD } (m, n, \sin, +, 0)$$

## Higher-Order Functions

$$\text{FOLD } (m, n, u, b, i) \stackrel{\triangle}{=} \textbf{if } m > n$$
$$\textbf{then } i$$
$$\textbf{else } b \; ((u \; m), (\text{FOLD } (m + 1, n, u, b, i)))$$

$$\sum_{i=m}^{n} \cos \; i : \text{FOLD } (m, n, \cos, +, 0)$$

## Higher-Order Functions

$$\text{FOLD } (m, n, u, b, i) \stackrel{\triangle}{=} \textbf{if } m > n$$
$$\textbf{then } i$$
$$\textbf{else } b \ ((u \ m), (\text{FOLD } (m + 1, n, u, b, i)))$$

$$\prod_{i=m}^{n} \sin \ i : \text{FOLD } (m, n, \sin, \times, 1)$$

# Lambda Expressions

Anonymous Functions

$$\sum_{i=m}^{n} i^2$$

# Lambda Expressions

Anonymous Functions

$$\sum_{i=m}^{n} i^2$$

$$\text{SQR } i \quad \stackrel{\triangle}{=} \quad i \times i$$

# Lambda Expressions

Anonymous Functions

$$\sum_{i=m}^{n} i^2 \;=\; \text{FOLD}\,(m, n, \text{SQR}, +, 0)$$

$$\text{SQR}\; i \;\overset{\triangle}{=}\; i \times i$$

# Lambda Expressions
Anonymous Functions

$$\sum_{i=m}^{n} i^2 \;=\; \text{FOLD}\,(m, n, (\lambda i\; i \times i), +, 0)$$

# Nesting, Free Variables, and Closures

$$(\lambda x \, 2 \times x) \, 3 \;\; = \;\; 6$$

# Nesting, Free Variables, and Closures

$$(\lambda x \, 2 \times x) \, 3 \;\; = \;\; 6$$

$$((\lambda x \, \lambda y \, x + y) \, 3) \, 4 \;\; = \;\; 7$$

# Nesting, Free Variables, and Closures

$$(\lambda x\ 2 \times x)\ 3 \quad = \quad 6$$

$$(\lambda x\ \lambda y\ x + y)\ 3 \quad = \quad ?$$

# Nesting, Free Variables, and Closures

$$(\lambda x \, 2 \times x) \, 3 \;\; = \;\; 6$$

$$(\lambda x \, \lambda y \, x + y) \, 3 \;\;\; = \;\;\; \langle \{x \mapsto 3\}, \lambda y \, x + y \rangle$$

*It is, of course, not excluded that the range of arguments or range of values of a function should consist wholly or partly of functions. The derivative, as this notion appears in the elementary differential calculus, is a familiar mathematical example of a function for which both ranges consist of functions.*

(p. 1 ¶4)

Church, A. (1941). *The Calculi of Lambda Conversion*. Princeton University Press, Princeton, NJ.

Gottfried Leibniz
|
Jacob Bernoulli
|
Johann Bernoulli
|
Leonhard Euler
|
Joseph Louis Lagrange
|
Simeon Poisson
|
Michel Chasles
|
Hubert Anson Newton
|
Eliakim Hastings Moore
|
Oswald Veblen
|
Alonzo Church

## Outline

## Derivatives

$$\frac{\mathrm{d}ax^2}{\mathrm{d}x} \rightsquigarrow 2ax$$

# Derivatives

$$\frac{\mathrm{d}ax^2}{\mathrm{d}x} \leadsto 2ax$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : \underbrace{f}_{\mathbb{R}\to\mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}\to\mathbb{R}}$$

## Derivatives

$$\frac{\mathrm{d}ax^2}{\mathrm{d}x} \rightsquigarrow 2ax$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : \underbrace{f}_{\mathbb{R}\to\mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}\to\mathbb{R}}$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : (\mathbb{R} \to \mathbb{R}) \to (\mathbb{R} \to \mathbb{R})$$

## Derivatives

$$\frac{\mathrm{d}ax^2}{\mathrm{d}x} \rightsquigarrow 2ax$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : \underbrace{f}_{\mathbb{R}\to\mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}\to\mathbb{R}}$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : (\mathbb{R} \to \mathbb{R}) \to (\mathbb{R} \to \mathbb{R})$$

$$\mathcal{D} : (\mathbb{R} \to \mathbb{R}) \to (\mathbb{R} \to \mathbb{R})$$

## Derivatives

$$\frac{\mathrm{d}ax^2}{\mathrm{d}x} \rightsquigarrow 2ax$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : \underbrace{f}_{\mathbb{R}\to\mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}\to\mathbb{R}}$$

$$\frac{\mathrm{d}}{\mathrm{d}x} : (\mathbb{R} \to \mathbb{R}) \to (\mathbb{R} \to \mathbb{R})$$

$$\mathcal{D} : (\mathbb{R} \to \mathbb{R}) \to (\mathbb{R} \to \mathbb{R})$$

$$\mathcal{D} \; \lambda x \; ax^2$$

## Partial Derivatives

$$\frac{\partial a x^2 y^3}{\partial x} \qquad\qquad \frac{\partial a x^2 y^3}{\partial y}$$

## Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x} \qquad\qquad \frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \; \lambda x \; ax^2y^3 \qquad\qquad \mathcal{D} \; \lambda y \; ax^2y^3$$

## Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x} \qquad\qquad \frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \; \lambda x \; ax^2y^3 \qquad\qquad \mathcal{D} \; \lambda y \; ax^2y^3$$

$$\mathcal{D}_1 \; \lambda(x,y) \; ax^2y^3 \qquad\qquad \mathcal{D}_2 \; \lambda(x,y) \; ax^2y^3$$

## Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x} \qquad\qquad \frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \; \lambda x \; ax^2y^3 \qquad\qquad \mathcal{D} \; \lambda y \; ax^2y^3$$

$$\mathcal{D}_1 \; \lambda(x,y) \; ax^2y^3 \qquad\qquad \mathcal{D}_2 \; \lambda(x,y) \; ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \to \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \to \mathbb{R}}$$

## Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x} \qquad\qquad \frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D} \; \lambda x \; ax^2y^3 \qquad\qquad \mathcal{D} \; \lambda y \; ax^2y^3$$

$$\mathcal{D}_1 \; \lambda(x,y) \; ax^2y^3 \qquad\qquad \mathcal{D}_2 \; \lambda(x,y) \; ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \to \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \to \mathbb{R}}$$

$$\frac{\partial}{\partial x} : (\mathbb{R}^n \to \mathbb{R}) \to (\mathbb{R}^n \to \mathbb{R})$$

## Partial Derivatives

$$\frac{\partial ax^2y^3}{\partial x} \qquad\qquad \frac{\partial ax^2y^3}{\partial y}$$

$$\mathcal{D}\ \lambda x\ ax^2y^3 \qquad\qquad \mathcal{D}\ \lambda y\ ax^2y^3$$

$$\mathcal{D}_1\ \lambda(x,y)\ ax^2y^3 \qquad\qquad \mathcal{D}_2\ \lambda(x,y)\ ax^2y^3$$

$$\frac{\partial}{\partial x} : \underbrace{f}_{\mathbb{R}^n \to \mathbb{R}} \mapsto \underbrace{f'}_{\mathbb{R}^n \to \mathbb{R}}$$

$$\frac{\partial}{\partial x} : (\mathbb{R}^n \to \mathbb{R}) \to (\mathbb{R}^n \to \mathbb{R})$$

$$\mathcal{D}_i : (\mathbb{R}^n \to \mathbb{R}) \to (\mathbb{R}^n \to \mathbb{R})$$

## Gradients

$$\nabla f \; \mathbf{x} \;\; = \;\; (\mathcal{D}_1 \, f \, \mathbf{x}), \ldots, (\mathcal{D}_n \, f \, \mathbf{x})$$

$$\nabla \;\; : \;\; (\mathbb{R}^n \to \mathbb{R}) \to (\mathbb{R}^n \to \mathbb{R}^n)$$

## Jacobians

$$f \quad : \quad \mathbb{R}^n \to \mathbb{R}^m$$

$$\mathbf{f} \quad : \quad (\mathbb{R}^n \to \mathbb{R})^m$$

$$(\mathcal{J} f \, \mathbf{x})[i,j] \quad = \quad (\nabla \, (\mathbf{f}[i]))[j]$$

$$\mathcal{J} \quad : \quad (\mathbb{R}^n \to \mathbb{R}^m) \to (\mathbb{R}^n \to \mathbb{R}^{m \times n})$$

## Operators

$\mathcal{D}$, $\nabla$, and $\mathcal{J}$ are traditionally called *operators*.
A more modern term is *higher-order functions*.
Higher-order functions are common in mathematics, physics, and engineering:

> *summations, comprehensions, quantifications, optimizations,*
> *integrals, convolutions, filters, edge detectors, Fourier transforms,*
> *differential equations, Hamiltonians, . . .*

# The Chain Rule

$$(f \circ g) \, x = g \, (f \, x)$$

## The Chain Rule

$$(f \circ g)\, x = g\,(f\,x)$$

$$\frac{\mathrm{d}g}{\mathrm{d}x} = \frac{\mathrm{d}g}{\mathrm{d}f}\frac{\mathrm{d}f}{\mathrm{d}x}$$

## The Chain Rule

$$(f \circ g)\, x = g\,(f\,x)$$

$$\frac{\mathrm{d}g}{\mathrm{d}x} = \frac{\mathrm{d}g}{\mathrm{d}f}\frac{\mathrm{d}f}{\mathrm{d}x}$$

$$\mathcal{D}\,(f \circ g)\, x = (\mathcal{D}\,g\,(f\,x)) \times (\mathcal{D}\,f\,x)$$

## The Chain Rule

$$(f \circ g)\ x = g\ (f\ x)$$

$$\frac{\mathrm{d}g}{\mathrm{d}x} = \frac{\mathrm{d}g}{\mathrm{d}f}\frac{\mathrm{d}f}{\mathrm{d}x}$$

$$\mathcal{D}\ (f \circ g)\ x = (\mathcal{D}\ g\ (f\ x)) \times (\mathcal{D}\ f\ x)$$

$$\mathcal{J}\ (f \circ g)\ \mathbf{x} = (\mathcal{J}\ g\ (f\ \mathbf{x})) \times (\mathcal{J}\ f\ \mathbf{x})$$

## Outline

1. Lambda Calculus

2. Differential Calculus in Lambda-Calculus Notation

3. Tutorial on AD
   - Forward Mode
   - Reverse Mode

4. Essence of the Derivation of Functional Reverse Mode

5. AD in Lambda-Calculus Notation

6. Examples

7. Benefits of this Approach

# Outline

1. Lambda Calculus

2. Differential Calculus in Lambda-Calculus Notation

3. Tutorial on AD
   - Forward Mode
   - Reverse Mode

4. Essence of the Derivation of Functional Reverse Mode

5. AD in Lambda-Calculus Notation

6. Examples

7. Benefits of this Approach

## Straight-Line Code and Jacobians

$$\mathbf{x}_1 = f_1 \; \mathbf{x}_0$$
$$\vdots$$
$$\mathbf{x}_n = f_n \; \mathbf{x}_{n-1}$$

$$f = f_1 \circ \cdots \circ f_n$$

$$\mathcal{J} f \; \mathbf{x}_0 = (\mathcal{J} f_n \; \mathbf{x}_{n-1}) \times \cdots \times (\mathcal{J} f_1 \; \mathbf{x}_0)$$
$$(\mathcal{J} f \; \mathbf{x}_0)^\top = (\mathcal{J} f_1 \; \mathbf{x}_0)^\top \times \cdots \times (\mathcal{J} f_n \; \mathbf{x}_{n-1})^\top$$

# One Way to Compute the Jacobian

$$\overline{\mathbf{X}_1'} = (\mathcal{J} f_1 \; \mathbf{x}_0)$$

$$\overline{\mathbf{X}_2'} = (\mathcal{J} f_2 \; \mathbf{x}_1) \times \overline{\mathbf{X}_1'}$$

$$\vdots$$

$$\overline{\mathbf{X}_n'} = (\mathcal{J} f_n \; \mathbf{x}_{n-1}) \times \overline{\mathbf{X}_{n-1}'}$$

$$\overline{\mathbf{X}_n'} = \mathcal{J} f \; \mathbf{x}_0$$

## Forward-Mode AD

$$\overline{\mathbf{x}_1'} = (\mathcal{J} f_1 \ \mathbf{x}_0) \times \overline{\mathbf{x}_0'}$$
$$\vdots$$
$$\overline{\mathbf{x}_n'} = (\mathcal{J} f_n \ \mathbf{x}_{n-1}) \times \overline{\mathbf{x}_{n-1}'}$$

$$\overline{\mathbf{x}_n'} = (\mathcal{J} f \ \mathbf{x}_0) \times \overline{\mathbf{x}_0'}$$

Wengert, R. E. (1964). A simple automatic derivative evaluation program. *Communications of the ACM*, **7**(8):463–4.

## Interleaving Forward Mode

$$\mathbf{x}_1 = f_1 \ \mathbf{x}_0 \qquad\qquad\qquad \overline{\mathbf{x}_1'} = (\mathcal{J} f_1 \ \mathbf{x}_0) \times \overline{\mathbf{x}_0'}$$

$$\vdots \qquad\qquad\qquad\qquad\qquad \vdots$$

$$\mathbf{x}_n = f_n \ \mathbf{x}_{n-1} \qquad\qquad\qquad \overline{\mathbf{x}_n'} = (\mathcal{J} f_n \ \mathbf{x}_{n-1}) \times \overline{\mathbf{x}_{n-1}'}$$

$$\mathbf{x}_1 = f_1 \ \mathbf{x}_0$$
$$\overline{\mathbf{x}_1'} = (\mathcal{J} f_1 \ \mathbf{x}_0) \times \overline{\mathbf{x}_0'}$$

$$\vdots$$

$$\mathbf{x}_n = f_n \ \mathbf{x}_{n-1}$$
$$\overline{\mathbf{x}_n'} = (\mathcal{J} f_n \ \mathbf{x}_{n-1}) \times \overline{\mathbf{x}_{n-1}'}$$

## Forward Mode as a Transformation

$$\left.\begin{array}{c} \mathbf{x}_1 = f_1\ \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n\ \mathbf{x}_{n-1} \end{array}\right\} \rightsquigarrow \left\{\begin{array}{c} \overrightarrow{\mathbf{x}_1} = \overrightarrow{f_1}\ \overrightarrow{\mathbf{x}_0} \\ \vdots \\ \overrightarrow{\mathbf{x}_n} = \overrightarrow{f_n}\ \overrightarrow{\mathbf{x}_{n-1}} \end{array}\right.$$
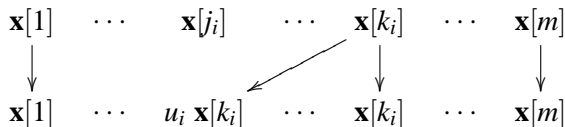
$$\overrightarrow{\mathbf{x}} = (\mathbf{x}, \overline{\mathbf{x}}')$$
$$\overrightarrow{f}\ (\mathbf{x}, \overline{\mathbf{x}}') = ((f\ \mathbf{x}), ((\mathcal{J}\ f\ \mathbf{x}) \times \overline{\mathbf{x}}'))$$
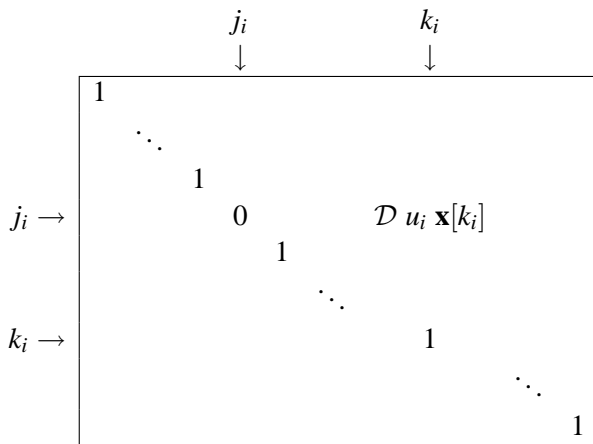
## A Unary Sparse Function

$$(f_i \, \mathbf{x})[j_i] = u_i \, \mathbf{x}[k_i]$$
$$(f_i \, \mathbf{x})[j'] = \mathbf{x}[j'] \qquad\qquad j' \neq j_i$$

$$\mathbf{x}[1] \quad \cdots \quad \mathbf{x}[j_i] \quad \cdots \quad \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[m]$$
$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$
$$\mathbf{x}[1] \quad \cdots \quad u_i \, \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[m]$$

# The Jacobian of a Unary Sparse Function

# Computing $(\mathcal{J} f_i \, \mathbf{x}_{i-1}) \times \overline{\mathbf{x}_{i-1}}$ for a Unary Sparse Function
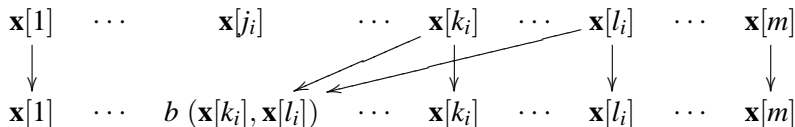
$$
\begin{pmatrix}
\overline{\mathbf{x}}[1] \\
\vdots \\
\overline{\mathbf{x}}[j_i - 1] \\
(\mathcal{D}\, u_i \, \mathbf{x}[k_i]) \times \overline{\mathbf{x}}[k_i] \\
\overline{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overline{\mathbf{x}}[k_i] \\
\vdots \\
\overline{\mathbf{x}}[m]
\end{pmatrix}
=
\begin{pmatrix}
1 \\
 & \ddots \\
 & & 1 \\
 & & & 0 & & \mathcal{D}\, u_i \, \mathbf{x}[k_i] \\
 & & & & 1 \\
 & & & & & \ddots \\
 & & & & & & 1 \\
 & & & & & & & \ddots \\
 & & & & & & & & 1
\end{pmatrix}
\begin{pmatrix}
\overline{\mathbf{x}}[1] \\
\vdots \\
\overline{\mathbf{x}}[j_i - 1] \\
\overline{\mathbf{x}}[j_i] \\
\overline{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overline{\mathbf{x}}[k_i] \\
\vdots \\
\overline{\mathbf{x}}[m]
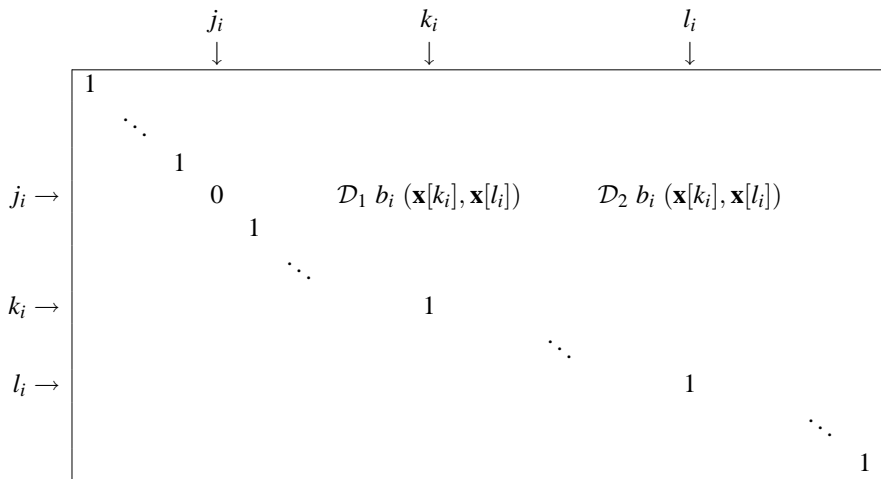\end{pmatrix}
$$

# A Binary Sparse Function

$$(f_i\ \mathbf{x})[j_i] = b_i\ (\mathbf{x}[k_i], \mathbf{x}[l_i])$$
$$(f_i\ \mathbf{x})[j'] = \mathbf{x}[j'] \qquad\qquad j' \neq j_i$$

# The Jacobian of a Binary Sparse Function

# Computing $\left(\mathcal{J} f_i \, \mathbf{x}_{i-1}\right) \times \overrightarrow{\mathbf{x}_{i-1}}$ for a Binary Sparse Function

$$
\begin{pmatrix}
\overline{\mathbf{x}}[1] \\
\vdots \\
\overline{\mathbf{x}}[j_i - 1] \\
((\mathcal{D}_1 \, b_i \, (\mathbf{x}[k_i], \mathbf{x}[l_i])) \times \overline{\mathbf{x}}[k_i]) + ((\mathcal{D}_2 \, b_i \, (\mathbf{x}[k_i], \mathbf{x}[l_i])) \times \overline{\mathbf{x}}[l_i]) \\
\overline{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overline{\mathbf{x}}[k_i] \\
\vdots \\
\overline{\mathbf{x}}[l_i] \\
\vdots \\
\overline{\mathbf{x}}[m]
\end{pmatrix}
=
\begin{pmatrix}
1 \\
 & \ddots \\
 & & 1 \\
 & & & 0 & & \mathcal{D}_1 \, b_i \, (\mathbf{x}[k_i], \mathbf{x}[l_i]) & & \mathcal{D}_2 \, b_i \, (\mathbf{x}[k_i], \mathbf{x}[l_i]) \\
 & & & & 1 \\
 & & & & & \ddots \\
 & & & & & & 1 \\
 & & & & & & & \ddots \\
 & & & & & & & & 1 \\
 & & & & & & & & & \ddots \\
 & & & & & & & & & & 1
\end{pmatrix}
\begin{pmatrix}
\overline{\mathbf{x}}[1] \\
\vdots \\
\overline{\mathbf{x}}[j_i - 1] \\
\overline{\mathbf{x}}[j_i] \\
\overline{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overline{\mathbf{x}}[k_i] \\
\vdots \\
\overline{\mathbf{x}}[l_i] \\
\vdots \\
\overline{\mathbf{x}}[m]
\end{pmatrix}
$$

## Forward Mode as a Sparse Transformation

$$x_{j_i} := u_i \, x_{k_i} \quad \rightsquigarrow \quad \overrightarrow{x_{j_i}} := \overrightarrow{u_i} \, \overrightarrow{x_{k_i}}$$
$$x_{j_i} := b_i \, (x_{k_i}, x_{l_i}) \quad \rightsquigarrow \quad \overrightarrow{x_{j_i}} := \overrightarrow{b_i} \, (\overrightarrow{x_{k_i}}, \overrightarrow{x_{l_i}})$$

$$\overrightarrow{x} = (x, \overrightarrow{x})$$
$$\overrightarrow{u} \, (x, \overrightarrow{x}) = ((u \, x), ((\mathcal{D} \, u \, x) \times \overrightarrow{x}))$$
$$\overrightarrow{b} \, ((x_1, \overrightarrow{x_1}), (x_2, \overrightarrow{x_2})) = ((b \, (x_1, x_2)),$$
$$(((\mathcal{D}_1 \, b \, (x_1, x_2)) \times \overrightarrow{x_1}) + ((\mathcal{D}_2 \, b \, (x_1, x_2)) \times \overrightarrow{x_2})))$$

# Outline

## Straight-Line Code and Jacobians

$$\mathbf{x}_1 = f_1 \ \mathbf{x}_0$$
$$\vdots$$
$$\mathbf{x}_n = f_n \ \mathbf{x}_{n-1}$$

$$f = f_1 \circ \cdots \circ f_n$$

$$\mathcal{J} f \ \mathbf{x}_0 = (\mathcal{J} f_n \ \mathbf{x}_{n-1}) \times \cdots \times (\mathcal{J} f_1 \ \mathbf{x}_0)$$
$$(\mathcal{J} f \ \mathbf{x}_0)^\top = (\mathcal{J} f_1 \ \mathbf{x}_0)^\top \times \cdots \times (\mathcal{J} f_n \ \mathbf{x}_{n-1})^\top$$

## Another Way to Compute the Jacobian

$$\overline{\mathbf{X}_{n-1}} = (\mathcal{J} f_n \, \mathbf{x}_{n-1})^\top$$
$$\overline{\mathbf{X}_{n-2}} = (\mathcal{J} f_{n-1} \, \mathbf{x}_{n-2})^\top \times \overline{\mathbf{X}_{n-1}}$$
$$\vdots$$
$$\overline{\mathbf{X}_0} = (\mathcal{J} f_1 \, \mathbf{x}_0)^\top \times \overline{\mathbf{X}_1}$$

$$\overline{\mathbf{X}_0} = (\mathcal{J} f \, \mathbf{x}_0)^\top$$

## Reverse-Mode AD

$$\overleftarrow{\mathbf{x}_{n-1}} = \left(\mathcal{J} f_n \; \mathbf{x}_{n-1}\right)^{\top} \times \overleftarrow{\mathbf{x}_n}$$

$$\vdots$$

$$\overleftarrow{\mathbf{x}_0} = \left(\mathcal{J} f_1 \; \mathbf{x}_0\right)^{\top} \times \overleftarrow{\mathbf{x}_1}$$

$$\overleftarrow{\mathbf{x}_0} = \left(\mathcal{J} f \; \mathbf{x}_0\right)^{\top} \times \overleftarrow{\mathbf{x}_n}$$

Speelpenning, B. (1980). *Compiling Fast Partial Derivatives of Functions Given by Algorithms*. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign.

# Reverse Mode Cannot be Interleaved

$$\mathbf{x}_1 = f_1 \, \mathbf{x}_0$$

$$\vdots$$

$$\mathbf{x}_n = f_n \, \mathbf{x}_{n-1}$$

$$\overleftarrow{\mathbf{x}_{n-1}} = \left(\mathcal{J} f_n \, \mathbf{x}_{n-1}\right)^\top \times \overleftarrow{\mathbf{x}_n}$$

$$\vdots$$

$$\overleftarrow{\mathbf{x}_0} = \left(\mathcal{J} f_1 \, \mathbf{x}_0\right)^\top \times \overleftarrow{\mathbf{x}_1}$$

# Reverse Mode via Backpropagators

$$\mathbf{x}_1 = f_1 \, \mathbf{x}_0$$
$$\overline{\mathbf{x}_1} = \lambda \overleftarrow{\mathbf{x}} \; \overline{\mathbf{x}_0} \; ((\mathcal{J} f_1 \, \mathbf{x}_0)^\top \times \overleftarrow{\mathbf{x}})$$
$$\vdots$$
$$\mathbf{x}_n = f_n \, \mathbf{x}_{n-1}$$
$$\overline{\mathbf{x}_n} = \lambda \overleftarrow{\mathbf{x}} \; \overline{\mathbf{x}_{n-1}} \; ((\mathcal{J} f_n \, \mathbf{x}_{n-1})^\top \times \overleftarrow{\mathbf{x}})$$

$$\overline{\mathbf{x}_n} \; \overleftarrow{\mathbf{x}_n}$$

# Reverse Mode as a Transformation

$$\left. \begin{array}{c} \mathbf{x}_1 = f_1 \ \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n \ \mathbf{x}_{n-1} \end{array} \right\} \rightsquigarrow \left\{ \begin{array}{c} \overleftarrow{\mathbf{x}_1} = \overleftarrow{f_1} \ \overleftarrow{\mathbf{x}_0} \\ \vdots \\ \overleftarrow{\mathbf{x}_n} = \overleftarrow{f_n} \ \overleftarrow{\mathbf{x}_{n-1}} \end{array} \right.$$

$$
\begin{aligned}
\overleftarrow{\mathbf{x}} &= (\mathbf{x}, \overline{\mathbf{x}}) \\
\overleftarrow{f} \ (\mathbf{x}, \overline{\mathbf{x}}) &= ((f \ \mathbf{x}), (\lambda \grave{\mathbf{x}} \ \overline{\mathbf{x}} \, ((\mathcal{J} f \ \mathbf{x})^\top \times \grave{\mathbf{x}})))
\end{aligned}
$$

# Reverse Mode via a Tape

$$\left. \begin{array}{c} \mathbf{x}_1 = f_1 \ \mathbf{x}_0 \\ \vdots \\ \mathbf{x}_n = f_n \ \mathbf{x}_{n-1} \end{array} \right\} \leadsto \left\{ \begin{array}{c} \overleftarrow{\mathbf{x}_1} = \overleftarrow{f_1} \ \overleftarrow{\mathbf{x}_0} \\ \vdots \\ \overleftarrow{\mathbf{x}_n} = \overleftarrow{f_n} \ \overleftarrow{\mathbf{x}_{n-1}} \end{array} \right.$$
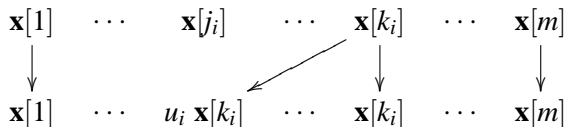
$$\begin{array}{rcl} \overleftarrow{\mathbf{x}} & = & \mathbf{x} \\ \overleftarrow{f} \ \mathbf{x} & = & \mathbf{begin} \ \bar{\mathbf{x}} := \lambda \overleftarrow{\mathbf{x}} \ \bar{\mathbf{x}} \left( (\mathcal{J} f \ \mathbf{x})^\top \times \overleftarrow{\mathbf{x}} \right); \\ & & \quad (f \ \mathbf{x}) \ \mathbf{end} \end{array}$$

# A Unary Sparse Function

$$(f_i \, \mathbf{x})[j_i] = u_i \, \mathbf{x}[k_i]$$
$$(f_i \, \mathbf{x})[j'] = \mathbf{x}[j'] \qquad\qquad j' \neq j_i$$



$\mathbf{x}[1] \quad \cdots \quad \mathbf{x}[j_i] \quad \cdots \quad \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[m]$

$\mathbf{x}[1] \quad \cdots \quad u_i \, \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[k_i] \quad \cdots \quad \mathbf{x}[m]$

# The Jacobian of a Unary Sparse Function

# The Transpose of the Jacobian of a Unary Sparse Function

$$
\begin{pmatrix}
1 & & & & & & & & \\
& \ddots & & & & & & & \\
& & 1 & & & & & & \\
& & & 0 & & \mathcal{D}\, u_i\, \mathbf{x}[k_i] & & & \\
& & & & 1 & & & & \\
& & & & & \ddots & & & \\
& & & & & & 1 & & \\
& & & & & & & \ddots & \\
& & & & & & & & 1
\end{pmatrix}^{\top}
=
\begin{pmatrix}
1 & & & & & & & & \\
& \ddots & & & & & & & \\
& & 1 & & & & & & \\
& & & 0 & & & & & \\
& & & & 1 & & & & \\
& & & & & \ddots & & & \\
& & \mathcal{D}\, u_i\, \mathbf{x}[k_i] & & 1 & & & \\
& & & & & & & \ddots & \\
& & & & & & & & 1
\end{pmatrix}
$$

# Computing $(\mathcal{J} f_i \, \mathbf{x}_{i-1})^\top \times \overleftarrow{\mathbf{x}_i}$ for a Unary Sparse Function
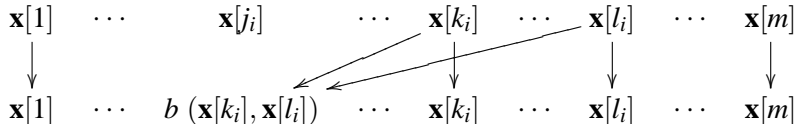
$$
\begin{pmatrix}
\overleftarrow{\mathbf{x}}[1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[j_i - 1] \\
0 \\
\overleftarrow{\mathbf{x}}[j_i + 1] \\
\vdots \\
((\mathcal{D}\, u_i\, \mathbf{x}[k_i]) \times \overleftarrow{\mathbf{x}}[j_i]) + \overleftarrow{\mathbf{x}}[k_i] \\
\vdots \\
\overleftarrow{\mathbf{x}}[m]
\end{pmatrix}
=
\begin{pmatrix}
1 \\
& \ddots \\
& & 1 \\
& & & 0 \\
& & & & 1 \\
& & & & & \ddots \\
& & & \mathcal{D}\, u_i\, \mathbf{x}[k_i] & & & 1 \\
& & & & & & & \ddots \\
& & & & & & & & 1
\end{pmatrix}
\begin{pmatrix}
\overleftarrow{\mathbf{x}}[1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[j_i - 1] \\
\overleftarrow{\mathbf{x}}[j_i] \\
\overleftarrow{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[k_i] \\
\vdots \\
\overleftarrow{\mathbf{x}}[m]
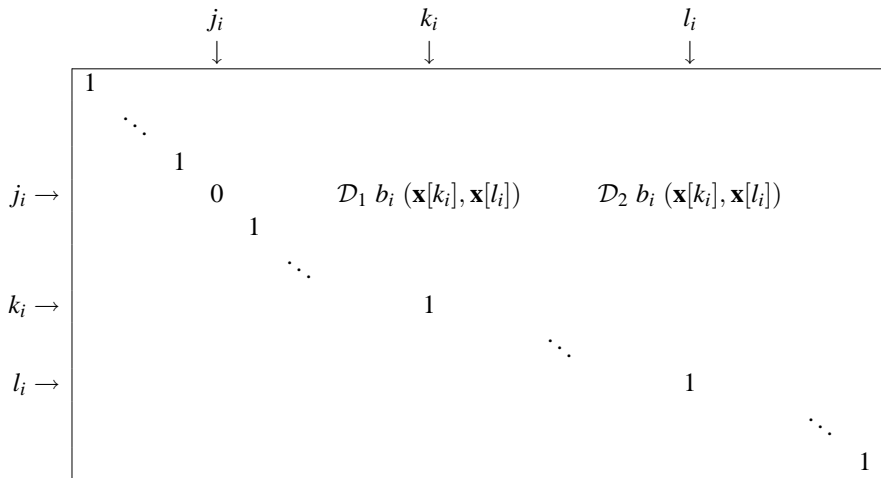\end{pmatrix}
$$

## A Binary Sparse Function

$$(f_i \ \mathbf{x})[j_i] = b_i \ (\mathbf{x}[k_i], \mathbf{x}[l_i])$$
$$(f_i \ \mathbf{x})[j'] = \mathbf{x}[j'] \qquad\qquad j' \neq j_i$$

## The Jacobian of a Binary Sparse Function

# The Transpose of the Jacobian of a Binary Sparse Function

$$
\begin{pmatrix}
1 & & & & & & & & & & \\
& \ddots & & & & & & & & & \\
& & 1 & & & & & & & & \\
& & & 0 & 1 & & \mathcal{D}_1\, b_i\, (\mathbf{x}[k_i], \mathbf{x}[l_i]) & & \mathcal{D}_2\, b_i\, (\mathbf{x}[k_i], \mathbf{x}[l_i]) & & \\
& & & & 1 & & & & & & \\
& & & & & \ddots & & & & & \\
& & & & & & 1 & & & & \\
& & & & & & & \ddots & & & \\
& & & & & & & & 1 & & \\
& & & & & & & & & \ddots & \\
& & & & & & & & & & 1
\end{pmatrix}^{\top}
=
\begin{pmatrix}
1 & & & & & & & & & & \\
& \ddots & & & & & & & & & \\
& & 1 & & & & & & & & \\
& & & 0 & & & & & & & \\
& & & & 1 & & & & & & \\
& & & & & \ddots & & & & & \\
& \mathcal{D}_1\, b_i\, (\mathbf{x}[k_i], \mathbf{x}[l_i]) & & & & & 1 & & & & \\
& & & & & & & \ddots & & & \\
& \mathcal{D}_2\, b_i\, (\mathbf{x}[k_i], \mathbf{x}[l_i]) & & & & & & & 1 & & \\
& & & & & & & & & \ddots & \\
& & & & & & & & & & 1
\end{pmatrix}
$$

# Computing $\left(\mathcal{J} f_i\, \mathbf{x}_{i-1}\right)^{\top} \times \overleftarrow{\mathbf{x}_i}$ for a Binary Sparse Function

$$
\begin{pmatrix}
\overleftarrow{\mathbf{x}}[1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[j_i - 1] \\
0 \\
\overleftarrow{\mathbf{x}}[j_i + 1] \\
\vdots \\
((\mathcal{D}_1\, b_i\,(\mathbf{x}[k_i], \mathbf{x}[l_i])) \times \overleftarrow{\mathbf{x}}[j_i]) + \overleftarrow{\mathbf{x}}[k_i] \\
\vdots \\
((\mathcal{D}_2\, b_i\,(\mathbf{x}[k_i], \mathbf{x}[l_i])) \times \overleftarrow{\mathbf{x}}[j_i]) + \overleftarrow{\mathbf{x}}[l_i] \\
\vdots \\
\overleftarrow{\mathbf{x}}[m]
\end{pmatrix}
=
\begin{pmatrix}
1 \\
 & \ddots \\
 & & 1 \\
 & & & 0 \\
 & & & & 1 \\
 & & & & & \ddots \\
 & & & \mathcal{D}_1\, b_i\,(\mathbf{x}[k_i], \mathbf{x}[l_i]) & & & 1 \\
 & & & & & & & \ddots \\
 & & & \mathcal{D}_2\, b_i\,(\mathbf{x}[k_i], \mathbf{x}[l_i]) & & & & & 1 \\
 & & & & & & & & & \ddots \\
 & & & & & & & & & & 1
\end{pmatrix}
\begin{pmatrix}
\overleftarrow{\mathbf{x}}[1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[j_i - 1] \\
\overleftarrow{\mathbf{x}}[j_i] \\
\overleftarrow{\mathbf{x}}[j_i + 1] \\
\vdots \\
\overleftarrow{\mathbf{x}}[k_i] \\
\vdots \\
\overleftarrow{\mathbf{x}}[l_i] \\
\vdots \\
\overleftarrow{\mathbf{x}}[m]
\end{pmatrix}
$$

## Sparse Reverse Mode via a Tape

$$x_{j_i} := u_i\ x_{k_i} \quad \rightsquigarrow \quad \bar{x} := \lambda[\,]\ \textbf{begin}\ \overline{x_{k_i}} +:= (\mathcal{D}\ u_i\ x_{k_i}) \times \overline{x_{j_i}};$$
$$\overline{x_{j_i}} := 0;$$
$$\bar{x}\ [\,]\ \textbf{end};$$
$$x_{j_i} := u_i\ x_{k_i}$$

$$x_{j_i} := b_i\ (x_{k_i}, x_{l_i}) \quad \rightsquigarrow \quad \bar{x} := \lambda[\,]\ \textbf{begin}\ \overline{x_{k_i}} +:= (\mathcal{D}_1\ b_i\ (x_{k_i}, x_{l_i})) \times \overline{x_{j_i}};$$
$$\overline{x_{l_i}} +:= (\mathcal{D}_2\ b_i\ (x_{k_i}, x_{l_i})) \times \overline{x_{j_i}};$$
$$\overline{x_{j_i}} := 0;$$
$$\bar{x}\ [\,]\ \textbf{end};$$
$$x_{j_i} := b_i\ (x_{k_i}, x_{l_i})$$

## Outline

# Reverse Mode on Imperative Programs

$$\vdots$$

$$x_2 \quad := \quad u\ x_1$$

$$x_4 \quad := \quad b\ (x_2, x_3)$$
$$\vdots$$

# Reverse Mode on Imperative Programs

$$\vdots$$

$$x_2 \quad := \quad u\ x_1$$

$$x_4 \quad := \quad b\ (x_2, x_3)$$

$$\vdots$$

$\left.\begin{array}{l} \\ \\ \\ \\ \\ \\ \end{array}\right\}$ *forward phase*

$$\vdots$$

$$\overleftarrow{x_2} \quad +:= \quad (\mathcal{D}_1\ b\ (x_2, x_3)) \times \overleftarrow{x_4}$$
$$\overleftarrow{x_3} \quad +:= \quad (\mathcal{D}_2\ b\ (x_2, x_3)) \times \overleftarrow{x_4}$$

$$\overleftarrow{x_1} \quad +:= \quad (\mathcal{D}\ u\ x_1) \times \overleftarrow{x_2}$$

$$\vdots$$

$\left.\begin{array}{l} \\ \\ \\ \\ \\ \\ \end{array}\right\}$ *reverse phase*

# Reverse Mode on Imperative Programs

$$\vdots$$

$$x_2 \quad := \quad u\ x_1$$

$$x_4 \quad := \quad b\ (x_2, x_3)$$
$$\vdots$$

$\left.\begin{array}{c} \\ \\ \\ \\ \\ \end{array}\right\}$ *forward phase*

$$\vdots$$

$$\overleftarrow{x_2} \quad +:= \quad (\mathcal{D}_1\ b\ (x_2, x_3)) \times \overleftarrow{x_4}$$
$$\overleftarrow{x_3} \quad +:= \quad (\mathcal{D}_2\ b\ (x_2, x_3)) \times \overleftarrow{x_4}$$

$$\overleftarrow{x_1} \quad +:= \quad (\mathcal{D}\ u\ x_1) \times \overleftarrow{x_2}$$
$$\vdots$$

$\left.\begin{array}{c} \\ \\ \\ \\ \\ \end{array}\right\}$ *reverse phase*

## Reverse Mode on Imperative Programs

$$
\left.
\begin{array}{rcl}
& \vdots & \\
x_2 & := & u\,x_1 \\
\\
x_4 & := & b\,(x_2, x_3) \\
& \vdots &
\end{array}
\right\}\ \textit{forward phase}
$$

$$
\left.
\begin{array}{rcl}
& \vdots & \\
\\
\overleftarrow{x_2} & \mathrel{+:=} & (\mathcal{D}_1\,b\,(x_2, x_3)) \times \overleftarrow{x_4} \\
\overleftarrow{x_3} & \mathrel{+:=} & (\mathcal{D}_2\,b\,(x_2, x_3)) \times \overleftarrow{x_4} \\
\\
\overleftarrow{x_1} & \mathrel{+:=} & (\mathcal{D}\,u\,x_1) \times \overleftarrow{x_2} \\
& \vdots &
\end{array}
\right\}\ \textit{reverse phase}
$$

# Reverse Mode on Imperative Programs

$$\left. \begin{array}{rcl} & \vdots & \\[1em] x_2 & := & u\ x_1 \\[2em] x_4 & := & b\ (x_2, x_3) \\ & \vdots & \end{array} \right\} \textit{forward phase}$$

$$\left. \begin{array}{rcl} & \vdots & \\[1em] \overleftarrow{x_2} & +:= & (\mathcal{D}_1\ b\ (x_2, x_3)) \times \overleftarrow{x_4} \\ \overleftarrow{x_3} & +:= & (\mathcal{D}_2\ b\ (x_2, x_3)) \times \overleftarrow{x_4} \\[1em] \overleftarrow{x_1} & +:= & (\mathcal{D}\ u\ x_1) \times \overleftarrow{x_2} \\ & \vdots & \end{array} \right\} \textit{reverse phase}$$

# Reverse Mode on Imperative Programs

$$
\begin{aligned}
&\qquad\qquad\vdots \\
x_2 \quad &:= \quad u\ x_1 \\
\\
x_2 \quad &:= \quad b\ (x_2, x_3) \\
&\qquad\qquad\vdots
\end{aligned}
\left.\begin{aligned}\\\\\\\\\\\\\\\end{aligned}\right\} \textit{forward phase}
$$

$$
\begin{aligned}
&\qquad\qquad\vdots \\
\\
\overleftarrow{x_2} \quad &+:= \quad (\mathcal{D}_1\ b\ (x_2, x_3)) \times \overleftarrow{x_2} \\
\overleftarrow{x_3} \quad &+:= \quad (\mathcal{D}_2\ b\ (x_2, x_3)) \times \overleftarrow{x_2} \\
\\
\overleftarrow{x_1} \quad &+:= \quad (\mathcal{D}\ u\ x_1) \times \overleftarrow{x_2} \\
&\qquad\qquad\vdots
\end{aligned}
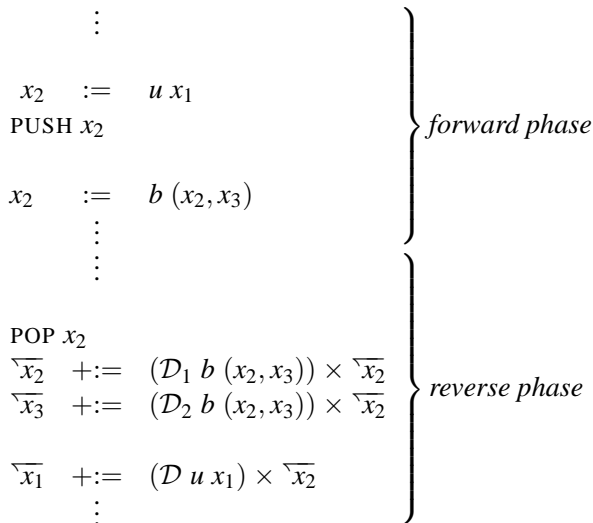\left.\begin{aligned}\\\\\\\\\\\\\\\end{aligned}\right\} \textit{reverse phase}
$$

# Reverse Mode on Imperative Programs

$$
\left.
\begin{array}{rcl}
 & \vdots & \\[1em]
x_2 & := & u \ x_1 \\[2em]
x_2 & := & b \ (x_2, x_3) \\
 & \vdots &
\end{array}
\right\} \textit{forward phase}
$$

$$
\left.
\begin{array}{rcl}
 & \vdots & \\[1em]
\overleftarrow{x_2} & +:= & (\mathcal{D}_1 \ b \ (x_2, x_3)) \times \overleftarrow{x_2} \\
\overleftarrow{x_3} & +:= & (\mathcal{D}_2 \ b \ (x_2, x_3)) \times \overleftarrow{x_2} \\[1em]
\overleftarrow{x_1} & +:= & (\mathcal{D} \ u \ x_1) \times \overleftarrow{x_2} \\
 & \vdots &
\end{array}
\right\} \textit{reverse phase}
$$

# Reverse Mode on Imperative Programs

$$
\left.
\begin{aligned}
&\qquad\vdots \\
x_2 \;\; &:= \;\; u\,x_1 \\
&\text{PUSH } x_2 \\
\\
x_2 \;\; &:= \;\; b\,(x_2,x_3) \\
&\qquad\vdots
\end{aligned}
\right\} \text{forward phase}
$$

$$
\left.
\begin{aligned}
&\qquad\vdots \\
&\text{POP } x_2 \\
\overleftarrow{x_2} \;\; &+:= \;\; (\mathcal{D}_1\,b\,(x_2,x_3)) \times \overleftarrow{x_2} \\
\overleftarrow{x_3} \;\; &+:= \;\; (\mathcal{D}_2\,b\,(x_2,x_3)) \times \overleftarrow{x_2} \\
\\
\overleftarrow{x_1} \;\; &+:= \;\; (\mathcal{D}\,u\,x_1) \times \overleftarrow{x_2} \\
&\qquad\vdots
\end{aligned}
\right\} \text{reverse phase}
$$

# Reverse Mode on Imperative Programs

$$
\left.
\begin{array}{l}
\quad\vdots \\
\text{PUSH } x_1 \\
\quad x_2 \quad := \quad u\ x_1 \\
\text{PUSH } x_2 \\
\text{PUSH } x_3 \\
x_2 \quad := \quad b\ (x_2, x_3) \\
\quad\vdots
\end{array}
\right\} \textit{forward phase}
$$

$$
\left.
\begin{array}{l}
\quad\vdots \\
\text{POP } x_3 \\
\text{POP } x_2 \\
\overleftarrow{x_2} \ +:= \ (\mathcal{D}_1\ b\ (x_2, x_3)) \times \overleftarrow{x_2} \\
\overleftarrow{x_3} \ +:= \ (\mathcal{D}_2\ b\ (x_2, x_3)) \times \overleftarrow{x_2} \\
\text{POP } x_1 \\
\overleftarrow{x_1} \ +:= \ (\mathcal{D}\ u\ x_1) \times \overleftarrow{x_2} \\
\quad\vdots
\end{array}
\right\} \textit{reverse phase}
$$

## Notation

In the following slides, I use **x**, **y**, **x**$_i$, and **y**$_i$ to denote tuples of scalar variables, i.e. $(x_{47}, x_{19}, x_{33})$.

## Notation

In the following slides, I use $\mathbf{x}$, $\mathbf{y}$, $\mathbf{x}_i$, and $\mathbf{y}_i$ to denote tuples of scalar variables, i.e. $(x_{47}, x_{19}, x_{33})$.

I use $\grave{\mathbf{x}}$, $\grave{\mathbf{y}}$, $\grave{\mathbf{x}}_i$, and $\grave{\mathbf{y}}_i$ to denote tuples of corresponding sensitivities of scalar variables, i.e. $(\grave{x_{47}}, \grave{x_{19}}, \grave{x_{33}})$.

# Subroutines and Tapes

Unary Primitives

$$u : x \mapsto y$$

# Subroutines and Tapes

Unary Primitives

$$u : x \mapsto y$$

$$\overleftarrow{u} : x \mapsto y \quad \triangleq \quad \begin{cases} \text{PUSH } x \\ y \;\; := \;\; u \, x \end{cases}$$

$$\overline{u} : \overleftarrow{y} \mapsto \overleftarrow{x} \quad \triangleq \quad \begin{cases} \text{POP } x \\ \overleftarrow{x} \;\; +\!:= \;\; (\mathcal{D} \, u \, x) \times \overleftarrow{y} \end{cases}$$

# Subroutines and Tapes

Binary Primitives

$$b : (x, y) \mapsto z$$

# Subroutines and Tapes
Binary Primitives

$$b : (x, y) \mapsto z$$

$$\overset{\leftarrow}{b} : (x, y) \mapsto z \quad \triangleq \quad \begin{cases} \text{PUSH } x \\ \text{PUSH } y \\ z \quad := \quad b\,(x, y) \end{cases}$$

$$\overline{b} : \overset{\backprime}{z} \mapsto (\overset{\backprime}{x}, \overset{\backprime}{y}) \quad \triangleq \quad \begin{cases} \text{POP } x \\ \text{POP } y \\ \overset{\backprime}{x} \quad +:= \quad (\mathcal{D}_1\, b\,(x, y)) \times \overset{\backprime}{z} \\ \overset{\backprime}{y} \quad +:= \quad (\mathcal{D}_2\, b\,(x, y)) \times \overset{\backprime}{z} \end{cases}$$

# Subroutines and Tapes

User-Defined Functions

$$f : \mathbf{x} \mapsto \mathbf{y} \qquad \triangleq \qquad \begin{cases} \mathbf{y}_1 & := & f_1 \ \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n & := & f_n \ \mathbf{x}_n \end{cases}$$

## Subroutines and Tapes
User-Defined Functions

$$f : \mathbf{x} \mapsto \mathbf{y} \quad \triangleq \quad \begin{cases} \mathbf{y}_1 & := & f_1 \ \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n & := & f_n \ \mathbf{x}_n \end{cases}$$

$$\overleftarrow{f} : \mathbf{x} \mapsto \mathbf{y} \quad \triangleq \quad \begin{cases} \mathbf{y}_1 & := & \overleftarrow{f_1} \ \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n & := & \overleftarrow{f_n} \ \mathbf{x}_n \end{cases}$$

$$\bar{f} : \overleftarrow{\mathbf{y}} \mapsto \overleftarrow{\mathbf{x}} \quad \triangleq \quad \begin{cases} \overline{\mathbf{x}_n} & +:= & \overline{f_n} \ \overleftarrow{\mathbf{y}_n} \\ & \vdots & \\ \overline{\mathbf{x}_1} & +:= & \overline{f_1} \ \overleftarrow{\mathbf{y}_1} \end{cases}$$

# Representing the Tape as Function Arguments and Results
Unary Primitives

$$u : x \mapsto y$$

$$\overleftarrow{u} : x \mapsto (y, x) \quad \triangleq \quad \{ \ y \ := \ u \, x$$

$$\overline{u} : (x, \overleftarrow{y}) \mapsto \overleftarrow{x} \quad \triangleq \quad \{ \ \overleftarrow{x} \ +:= \ (\mathcal{D} \, u \, x) \times \overleftarrow{y}$$

# Representing the Tape as Function Arguments and Results
Binary Primitives

$$b : (x, y) \mapsto z$$

$$\overleftarrow{b} : (x, y) \mapsto (z, (x, y)) \quad \triangleq \quad \{\ z\ :=\ b\,(x, y)$$

$$\overline{b} : ((x, y), \overleftarrow{z}) \mapsto (\overleftarrow{x}, \overleftarrow{y}) \quad \triangleq \quad \left\{ \begin{array}{rl} \overleftarrow{x} & +:= \ (\mathcal{D}_1\,b\,(x, y)) \times \overleftarrow{z} \\ \overleftarrow{y} & +:= \ (\mathcal{D}_2\,b\,(x, y)) \times \overleftarrow{z} \end{array} \right.$$

# Representing the Tape as Function Arguments and Results
## User-Defined Functions

$$f : \mathbf{x} \mapsto \mathbf{y} \qquad \triangleq \qquad \left\{ \begin{array}{rcl} \mathbf{y}_1 & := & f_1 \; \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n & := & f_n \; \mathbf{x}_n \end{array} \right.$$

$$\overleftarrow{f} : \mathbf{x} \mapsto (\mathbf{y}, (\mathbf{t}_1, \ldots, \mathbf{t}_n)) \qquad \triangleq \qquad \left\{ \begin{array}{rcl} \mathbf{y}_1, \mathbf{t}_1 & := & \overleftarrow{f_1} \; \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n, \mathbf{t}_n & := & \overleftarrow{f_n} \; \mathbf{x}_n \end{array} \right.$$

$$\bar{f} : ((\mathbf{t}_1, \ldots, \mathbf{t}_n), \overleftarrow{\mathbf{y}}) \mapsto \overleftarrow{\mathbf{x}} \qquad \triangleq \qquad \left\{ \begin{array}{rcl} \overleftarrow{\mathbf{x}_n} & \mathrel{+:=} & \overline{f_n} \; (\mathbf{t}_n, \overleftarrow{\mathbf{y}_n}) \\ & \vdots & \\ \overleftarrow{\mathbf{x}_1} & \mathrel{+:=} & \overline{f_1} \; (\mathbf{t}_1, \overleftarrow{\mathbf{y}_1}) \end{array} \right.$$

# Representing the Tape as Closures
Unary Primitives

$$u : x \mapsto y$$

$$\overleftarrow{u} : x \mapsto (y, \overline{u}) \quad \triangleq \quad \left\{ \begin{array}{lll} y & := & u\, x \\ \overline{u} : \overleftarrow{y} \mapsto \overleftarrow{x} & \triangleq & \{\ \overleftarrow{x} \ +:= \ (\mathcal{D}\, u\, x) \times \overleftarrow{y} \end{array} \right.$$

# Representing the Tape as Closures

Binary Primitives

$$b : (x, y) \mapsto z$$

$$\overleftarrow{b} : (x, y) \mapsto (z, \overline{b}) \quad \triangleq \quad \begin{cases} z & := b\,(x, y) \\ \overline{b} : \overleftarrow{z} \mapsto (\overleftarrow{x}, \overleftarrow{y}) \triangleq \begin{cases} \overleftarrow{x} & += (\mathcal{D}_1\, b\,(x, y)) \times \overleftarrow{z} \\ \overleftarrow{y} & += (\mathcal{D}_2\, b\,(x, y)) \times \overleftarrow{z} \end{cases} \end{cases}$$

# Representing the Tape as Closures
## User-Defined Functions

$$f : \mathbf{x} \mapsto \mathbf{y} \quad\triangleq\quad \begin{cases} \mathbf{y}_1 & := & f_1 \; \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n & := & f_n \; \mathbf{x}_n \end{cases}$$

$$\overleftarrow{f} : \mathbf{x} \mapsto (\mathbf{y}, \bar{f}) \quad\triangleq\quad \begin{cases} \mathbf{y}_1, \overline{f_1} & := & \overleftarrow{f_1} \; \mathbf{x}_1 \\ & \vdots & \\ \mathbf{y}_n, \overline{f_n} & := & \overleftarrow{f_n} \; \mathbf{x}_n \\ \bar{f} : \overleftarrow{\mathbf{y}} \mapsto \overleftarrow{\mathbf{x}} & \triangleq & \begin{cases} \overleftarrow{\mathbf{x}_n} & +:= & \overline{f_n} \; \overleftarrow{\mathbf{y}_n} \\ & \vdots & \\ \overleftarrow{\mathbf{x}_1} & +:= & \overline{f_1} \; \overleftarrow{\mathbf{y}_1} \end{cases} \end{cases}$$

# Details for Handling Closures Omitted

# Outline

# Traditional Formulation of AD as Transformations

Forward Mode: $\mathbb{R}^n \to \mathbb{R}^m \leadsto (\mathbb{R}^n \times \mathbb{R}^n) \to (\mathbb{R}^m \times \mathbb{R}^m)$

Reverse Mode: $\mathbb{R}^n \to \mathbb{R}^m \leadsto (\mathbb{R}^n \to (\mathbb{R}^m \times \mathbb{R}^l)) \times ((\mathbb{R}^m \times \mathbb{R}^l) \to \mathbb{R}^n)$

# New Formulation of AD as Higher-Order Functions
## Perturbation Types

$$\overline{\mathbf{null}} = \mathbf{null}$$

$$\overline{\mathbb{R}} = \mathbb{R}$$

$$\overline{\tau_1 \times \tau_2} = \overline{\tau_1'} \times \overline{\tau_2'}$$

$$\overline{\tau_1 \xrightarrow{\tau_1', \ldots, \tau_n'} \tau_2} = \overline{\tau_1'} \times \cdots \times \overline{\tau_n'}$$

# New Formulation of AD as Higher-Order Functions
Forward Types

$$\overrightarrow{\mathbf{null}} \;=\; \mathbf{null} \times \overrightarrow{\mathbf{null}}$$

$$\overrightarrow{\mathbb{R}} \;=\; \mathbb{R} \times \overrightarrow{\mathbb{R}}$$

$$\overrightarrow{\tau_1 \times \tau_2} \;=\; \overrightarrow{\tau_1} \times \overrightarrow{\tau_2}$$

$$\overrightarrow{\tau_1 \xrightarrow{\tau_1', \ldots, \tau_n'} \tau_2} \;=\; \overrightarrow{\tau_1} \xrightarrow{\overrightarrow{\tau_1'}, \ldots, \overrightarrow{\tau_n'}} \overrightarrow{\tau_2}$$

# New Formulation of AD as Higher-Order Functions
Sensitivity Types

$$\overleftarrow{\mathbf{null}} = \mathbf{null}$$

$$\overleftarrow{\mathbb{R}} = \mathbb{R}$$

$$\overleftarrow{\tau_1 \times \tau_2} = \overleftarrow{\tau_1} \times \overleftarrow{\tau_2}$$

$$\overleftarrow{\tau_1 \xrightarrow{\tau_1', \dots, \tau_n'} \tau_2} = \overleftarrow{\tau_1'} \times \cdots \times \overleftarrow{\tau_n'}$$

# New Formulation of AD as Higher-Order Functions
Reverse Types

$$\overleftarrow{\mathbf{null}} = \mathbf{null}$$

$$\overleftarrow{\mathbb{R}} = \mathbb{R}$$

$$\overleftarrow{\tau_1 \times \tau_2} = \overleftarrow{\tau_1} \times \overleftarrow{\tau_2}$$

$$\overleftarrow{\tau_1 \overset{\tau_1', \ldots, \tau_n'}{\longrightarrow} \tau_2} = \overleftarrow{\tau_1} \overset{\overleftarrow{\tau_1'}, \ldots, \overleftarrow{\tau_n'}}{\longrightarrow} (\overleftarrow{\tau_2} \times (\overleftarrow{\tau_2} \to (\overleftarrow{\tau_1'} \times \cdots \times \overleftarrow{\tau_n'}) \times \overleftarrow{\tau_1}))$$

# New Formulation of AD as Higher-Order Functions

Forward Mode: $\overrightarrow{\mathcal{J}} : \tau \to \overrightarrow{\tau}$

Reverse Mode: $\overleftarrow{\mathcal{J}} : \tau \to \overleftarrow{\tau}$

# Outline

## Derivatives

$$\mathcal{D} f \, x \;\; \triangleq \;\; \text{TANGENT} \, ((\overrightarrow{\mathcal{J}} f) \, (x \blacktriangleright 1))$$

$$\mathcal{D} f \, x \;\; \triangleq \;\; \text{CDR} \, ((\text{CDR} \, ((\overleftarrow{\mathcal{J}} f) \, (\overleftarrow{\mathcal{J}} x))) \, 1)$$

# Roots using Newton-Raphson

$$\text{ROOT } (f, x_0, \epsilon) \triangleq \textbf{let } x' \triangleq x_0 - \frac{f\, x_0}{\mathcal{D} f\, x_0}$$
$$\textbf{in if } |x_0 - x'| \leq \epsilon \textbf{ then } x_0 \textbf{ else } \text{ROOT } (f, x', \epsilon)$$

# Univariate Minimizer
Line Search

$$\text{LINESEARCH} \; (f, x_0, \epsilon) \stackrel{\triangle}{=} \text{ROOT} \; ((\mathcal{D} f), x_0, \epsilon)$$

# Gradients

$$\nabla f\, x \;\overset{\triangle}{=}\; \textbf{let } n \overset{\triangle}{=} \text{LENGTH } x$$
$$\textbf{in } \text{MAP } ((\lambda i\; \text{TANGENT } ((\overrightarrow{\mathcal{J}} f)\; (x \blacktriangleright e_{i,n}))), (\iota\, n))$$
$$\nabla f\, x \;\overset{\triangle}{=}\; \text{CDR } ((\text{CDR } ((\overleftarrow{\mathcal{J}} f)\; (\overleftarrow{\mathcal{J}} x)))\; 1)$$

# Multivariate Minimizer
Gradient Descent

$\text{GRADIENTDESCENT } (f, x_0, \epsilon) \triangleq$
  $\textbf{let } g \triangleq \nabla f \, x_0$
  $\textbf{in if } \|g\| \leq \epsilon$
    $\textbf{then } x_0$
    $\textbf{else GRADIENTDESCENT}$
      $(f, (x_0 + ((\text{LINESEARCH } ((\lambda k \, f \, (x_0 + (k \times g))), \epsilon)) \times g)), \epsilon)$

# Saddle Points
Continuous Two-Person Zero Sum Games

$\mathbf{x} : \mathbb{R}^m$

$\mathbf{y} : \mathbb{R}^n$

$\text{PAYOFF} : \mathbb{R}^m \times \mathbb{R}^n \to \mathbb{R}$

$\min_{\mathbf{x}} \max_{\mathbf{y}} \text{PAYOFF} (\mathbf{x}, \mathbf{y})$

# Saddle Points
Continuous Two-Person Zero Sum Games

$\mathbf{x} : \mathbb{R}^m$

$\mathbf{y} : \mathbb{R}^n$

$\text{PAYOFF} : \mathbb{R}^m \times \mathbb{R}^n \to \mathbb{R}$

$\min_{\mathbf{x}} \max_{\mathbf{y}} \text{PAYOFF}\,(\mathbf{x}, \mathbf{y})$

$$(\mathbf{x}^*, \mathbf{y}^*) = \textbf{let } \mathbf{x}^* \triangleq \text{ARGMIN}\,((\lambda\mathbf{x}\,\text{MAX}\,((\lambda\mathbf{y}\,\text{PAYOFF}\,(\mathbf{x}, \mathbf{y})), \mathbf{y}_0, \epsilon)), \mathbf{x}_0, \epsilon)$$
$$\textbf{in } (\mathbf{x}^*, (\text{ARGMAX}\,((\lambda\mathbf{y}\,\text{PAYOFF}\,(\mathbf{x}^*, \mathbf{y})), \mathbf{y}_0, \epsilon)))$$

# Saddle Points
Continuous Two-Person Zero Sum Games
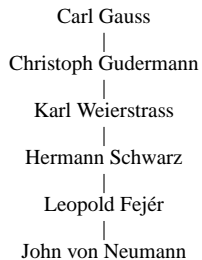
$\mathbf{x} : \mathbb{R}^m$

$\mathbf{y} : \mathbb{R}^n$

$\text{PAYOFF} : \mathbb{R}^m \times \mathbb{R}^n \to \mathbb{R}$

$\min\limits_{\mathbf{x}} \max\limits_{\mathbf{y}} \text{PAYOFF}(\mathbf{x}, \mathbf{y})$

$$(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{let}\ \mathbf{x}^* \stackrel{\triangle}{=} \text{ARGMIN}((\lambda\mathbf{x}\ \text{MAX}((\lambda\mathbf{y}\ \text{PAYOFF}(\mathbf{x}, \mathbf{y})), \mathbf{y}_0, \epsilon)), \mathbf{x}_0, \epsilon)$$
$$\mathbf{in}\ (\mathbf{x}^*, (\text{ARGMAX}((\lambda\mathbf{y}\ \text{PAYOFF}(\mathbf{x}^*, \mathbf{y})), \mathbf{y}_0, \epsilon)))$$

von Neumann, J. and Morgenstern, O. (1944). *Theory of Games and Economic Behavior.* Princeton University Press, Princeton, NJ.

Carl Gauss
|
Christoph Gudermann
|
Karl Weierstrass
|
Hermann Schwarz
|
Leopold Fejér
|
John von Neumann

# Function Inversion

$$f^{-1}\, y \overset{\triangle}{=} \text{ROOT}\left((\lambda x\, |(f\, x) - y|), x_0, \epsilon\right)$$

## Neural Nets

$$\text{NEURON } (\mathbf{w}, \mathbf{x}) \stackrel{\triangle}{=} \text{SIGMOID } (\mathbf{w} \cdot \mathbf{x})$$

$$\text{NEURALNET } ([\mathbf{w}''; \mathbf{w}_1'; \ldots; \mathbf{w}_m'], \mathbf{x}) \stackrel{\triangle}{=}$$
$$\quad \text{NEURON } (\mathbf{w}'', [\text{NEURON } (\mathbf{w}_1', \mathbf{x}); \ldots; \text{NEURON } (\mathbf{w}_m', \mathbf{x})])$$

$$\text{ERROR } \mathbf{w} \stackrel{\triangle}{=}$$
$$\quad \|[y_1; \ldots; y_n] - [\text{NEURALNET } (\mathbf{w}, \mathbf{x}_1); \ldots; \text{NEURALNET } (\mathbf{w}, \mathbf{x}_n)]\|$$

$$\text{GRADIENTDESCENT } (\text{ERROR}, \mathbf{w}_0, \epsilon)$$

Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, **323**:533–6.

# Supervised Machine Learning
Function Approximation

$$\text{ERROR } \theta \stackrel{\triangle}{=} \|[y_1; \ldots; y_n] - [f(\theta, \mathbf{x}_1); \ldots; f(\theta, \mathbf{x}_n)]\|$$

$$\text{GRADIENTDESCENT } (\text{ERROR}, \theta_0, \epsilon)$$

# Maximum Likelihood Estimation

$$\text{GRADIENTDESCENT}\left(\left(\lambda\theta\left(-\prod_{\mathbf{x}\in\mathcal{X}}P(\mathbf{x}|\theta)\right)\right),\theta_0,\epsilon\right)$$

Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philos. Trans. Roy. Soc. London Ser. A*, **222**:309–68.

# Engineering Design

PERFORMANCEOF SPLINECONTROLPOINTS $\triangleq$
  **let** WING $\triangleq$ SPLINETOSURFACE SPLINECONTROLPOINTS;
    AIRFLOW $\triangleq$ PDEsolver (WING, NAVIERSTOKES);
    LIFT, DRAG $\triangleq$ SURFACEINTEGRAL (WING, AIRFLOW, FORCE);
    PERFORMANCE $\triangleq$ DESIGNMETRIC (LIFT, DRAG, (WEIGHT WING))
  **in** PERFORMANCE

GRADIENTDESCENT (PERFORMANCEOF, SPLINECONTROLPOINTS$_0$, $\epsilon$)

# Outline

1. Lambda Calculus

2. Differential Calculus in Lambda-Calculus Notation

3. Tutorial on AD
   - Forward Mode
   - Reverse Mode

4. Essence of the Derivation of Functional Reverse Mode

5. AD in Lambda-Calculus Notation

6. Examples

7. Benefits of this Approach

# Our Approach is Efficient

# Our Approach is Efficient

- source-to-source transformation

# Our Approach is Efficient

- source-to-source transformation
- no overloading
- no interpretation of tape

# Our Approach is Efficient

- source-to-source transformation
- no overloading
- no interpretation of tape
- transformation conceptually done reflectively at run-time

## Our Approach is Efficient

- source-to-source transformation
- no overloading
- no interpretation of tape
- transformation conceptually done reflectively at run-time
- sophisticated compilation techniques can move transformation to compile-time

# Our Approach Exhibits Closure Properties

# Our Approach Exhibits Closure Properties

- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to *any* function

# Our Approach Exhibits Closure Properties

- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to *any* function including $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ themselves.

# Our Approach Exhibits Closure Properties

- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to *any* function including $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ themselves.
- The output of $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ are functions.

# Our Approach Exhibits Closure Properties

- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to *any* function
  including $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ themselves.
- The output of $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ are functions.
- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to the output of $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$.

# Our Approach Exhibits Closure Properties

- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to *any* function
  including $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ themselves.
- The output of $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ are functions.
- Can apply $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$ to the output of $\overrightarrow{\mathcal{J}}$ and $\overleftarrow{\mathcal{J}}$.
- Can take derivatives of arbitrary order.

# Our Approach Enhances Modularity

$$\text{ARGMIN}_1 \; (f, f') \quad \stackrel{\triangle}{=} \quad \ldots$$

# Our Approach Enhances Modularity

$$\textsc{Argmin}_1 \ (f, f') \quad \triangleq \quad \ldots$$

$$\ldots \textsc{Argmin}_1 \ (f, f') \ldots$$

# Our Approach Enhances Modularity

$$\text{ARGMIN}_1\ (f, f') \quad \triangleq \quad \dots$$
$$\text{ARGMIN}_2\ (f, f', f'') \quad \triangleq \quad \dots$$

$$\dots \text{ARGMIN}_1\ (f, f') \dots$$

# Our Approach Enhances Modularity

$$\mathrm{ARGMIN}_1\ (f, f') \quad \triangleq \quad \ldots$$
$$\mathrm{ARGMIN}_2\ (f, f', f'') \quad \triangleq \quad \ldots$$

$$\ldots \mathrm{ARGMIN}_2\ (f, f', f'') \ldots$$

# Our Approach Enhances Modularity

$$\text{ARGMIN}_1 f \quad \stackrel{\triangle}{=} \quad \ldots (\overleftrightarrow{\mathcal{J}} f) \ldots$$

# Our Approach Enhances Modularity

$$\text{ARGMIN}_1 f \;\; \triangleq \;\; \ldots (\overleftrightarrow{\mathcal{J}} f) \ldots$$

$$\ldots \text{ARGMIN}_1 f \ldots$$

# Our Approach Enhances Modularity

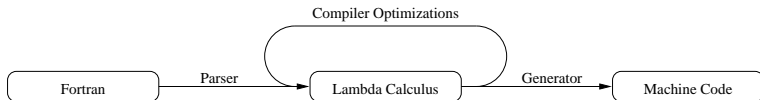$$\text{ARGMIN}_1 f \;\; \triangleq \;\; \dots (\overleftrightarrow{\mathcal{J}} f) \dots$$

$$\text{ARGMIN}_2 f \;\; \triangleq \;\; \dots (\overleftrightarrow{\mathcal{J}} f) \dots (\overleftrightarrow{\mathcal{J}} (\overleftrightarrow{\mathcal{J}} f)) \dots$$

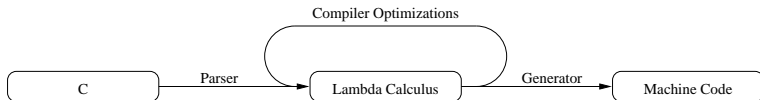$$\dots \text{ARGMIN}_1 f \dots$$

# Our Approach Enhances Modularity

$$\mathrm{ARGMIN}_1 \, f \;\; \triangleq \;\; \ldots (\overleftrightarrow{\mathcal{J}} \, f) \ldots$$

$$\mathrm{ARGMIN}_2 \, f \;\; \triangleq \;\; \ldots (\overleftrightarrow{\mathcal{J}} \, f) \ldots (\overleftrightarrow{\mathcal{J}} \, (\overrightarrow{\mathcal{J}} \, f)) \ldots$$
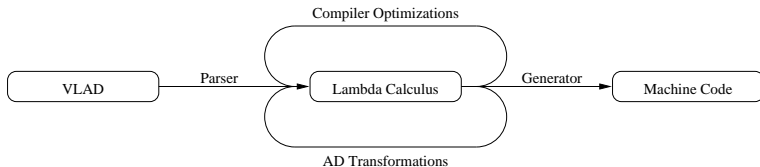
$$\ldots \mathrm{ARGMIN}_2 \, f \ldots$$

# Lambda the Ultimate Intermediate Language

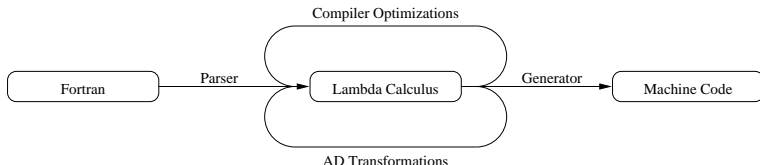# Lambda the Ultimate Intermediate Language

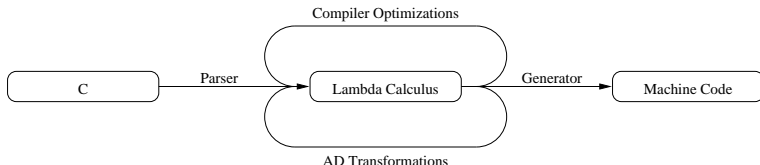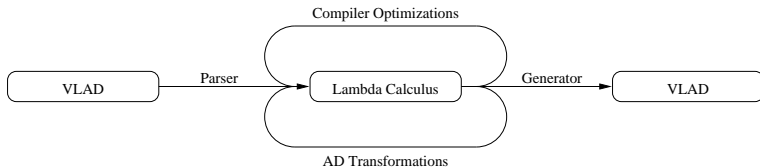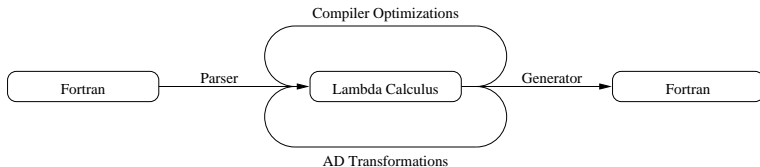# Lambda the Ultimate Intermediate Language

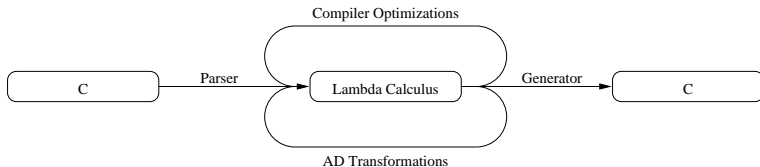# Lambda the Ultimate Intermediate Language *for AD*

# Lambda the Ultimate Intermediate Language *for AD*

# Lambda the Ultimate Intermediate Language *for AD*

# Lambda the Ultimate Intermediate Language *for AD*

# Lambda the Ultimate Intermediate Language *for AD*

# Lambda the Ultimate Intermediate Language *for AD*

# Our Work

# Our Work

- Prior Work

# Our Work

- Prior Work

  STALIN compiler for SCHEME

# Our Work

- Prior Work

  STALIN compiler for SCHEME
  ruthless, brutal, good at execution

# Our Work

- Prior Work

    STALIN compiler for SCHEME
    ruthless, brutal, good at execution
    $20 \times$ FORTRAN

## Our Work

- Prior Work

  STALIN compiler for SCHEME
  ruthless, brutal, good at execution
  $20 \times$ FORTRAN

- Current Work

# Our Work

- Prior Work

  STALIN compiler for SCHEME
  ruthless, brutal, good at execution
  $20 \times$ FORTRAN

- Current Work

  theory: $\lambda\nabla$-calculus

## Our Work

- Prior Work

  STALIN    compiler for SCHEME

           ruthless, brutal, good at execution

           $20 \times$ FORTRAN

- Current Work

  theory:   $\lambda\nabla$-calculus

         $\lambda$-calculus $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$

# Our Work

- Prior Work

  STALIN  compiler for SCHEME
          ruthless, brutal, good at execution
          $20 \times$ FORTRAN

- Current Work

  theory:  $\lambda\nabla$-calculus
         $\lambda$-calculus $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
  language:  VLAD

## Our Work

- Prior Work

    STALIN compiler for SCHEME
    ruthless, brutal, good at execution
    $20 \times$ FORTRAN

- Current Work

    theory: $\lambda\nabla$-calculus
    $\lambda\text{-calculus} + \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$

    language: VLAD
    $\text{SCHEME} + \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$

# Our Work

- Prior Work

  STALIN  compiler for SCHEME
           ruthless, brutal, good at execution
           $20 \times$ FORTRAN

- Current Work

  theory:  $\lambda\nabla$-calculus
         $\lambda$-calculus $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
  language:  VLAD
         SCHEME $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
         <u>F</u>unctional <u>L</u>anguage for <u>AD</u>

## Our Work

- Prior Work

  STALIN compiler for SCHEME
  ruthless, brutal, good at execution
  $20 \times$ FORTRAN

- Current Work

  theory: $\lambda\nabla$-calculus
  $\lambda$-calculus $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$

  language: VLAD
  SCHEME $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
  <u>F</u>unctional <u>L</u>anguage for <u>AD</u>

  compiler: STALIN$\nabla$

## Our Work

- Prior Work

    STALIN compiler for SCHEME
    
    ruthless, brutal, good at execution
    
    $20 \times$ FORTRAN

- Current Work

    theory: $\lambda\nabla$-calculus
    
    $\lambda$-calculus $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
    
    language: VLAD
    
    SCHEME $+ \overrightarrow{\mathcal{J}} + \overleftarrow{\mathcal{J}}$
    
    Functional Language for AD
    
    compiler: STALIN$\nabla$

manuscripts and code:

http://www-bcl.cs.nuim.ie/~qobi/stalingrad/