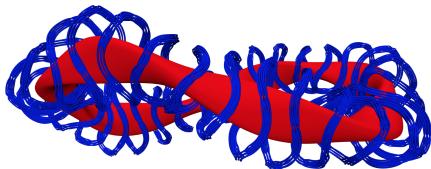


# Finite-Build Stellarator Coil Design & Automatic Differentiation

Nick McGreivy<sup>1,2</sup>   Stuart Hudson<sup>1</sup>   Caoxiang Zhu<sup>1</sup>



<sup>1</sup>Princeton Plasma Physics Laboratory

<sup>2</sup>PhD Candidate  
Program in Plasma Physics  
Princeton University

July 24th 2020

# Automatic Differentiation (AD)

Also known as algorithmic differentiation or computational differentiation

- Automatic Differentiation (AD) is a technology for automatically computing the exact numerical derivatives of any differentiable function  $\mathbf{y} = \mathbf{f}(\mathbf{x})$ , including arbitrarily complex simulations, represented by a computer program.
- To compute automatic derivatives of a function, you need to program that function using an AD software tool.

## Part I: Finite-build coil design & results

- FOCUSADD is a new finite-build coil design code using AD.
- Finite-build can be important when coil tolerances are small and/or when coils are large.

## Part II: Discussion of AD

- AD can be understood as (i) Compositions of primitive operations and (ii) Jacobian products.
- AD can help us better understand the adjoint method, and vice versa.

# Part I: Finite-build coil design & results



# Stellarator Coil Design as an Optimization Problem

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{p})$$

$$f(\mathbf{p}) = \int_S (\mathbf{B}(\mathbf{p}) \cdot \hat{\mathbf{n}})^2 dA + \text{Regularization term(s)}$$

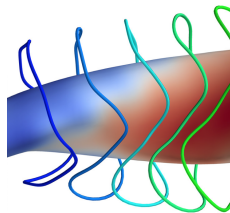
$$\mathbf{B} = \mathbf{B}_{\text{coils}} + \mathbf{B}_{\text{plasma}} \rightarrow 0$$

# FOCUS: Coil Design Code

$$x^i(\theta) = \sum_{m=0}^{N_F-1} X_{cm}^i \cos(m\theta) + X_{sm}^i \sin(m\theta)$$

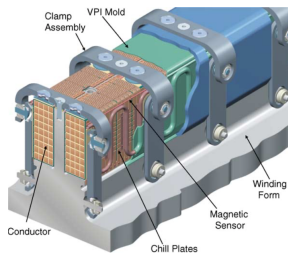
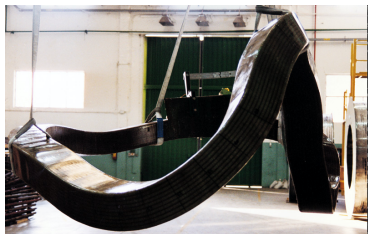
$$y^i(\theta) = \sum_{m=0}^{N_F-1} Y_{cm}^i \cos(m\theta) + Y_{sm}^i \sin(m\theta)$$

$$z^i(\theta) = \sum_{m=0}^{N_F-1} Z_{cm}^i \cos(m\theta) + Z_{sm}^i \sin(m\theta)$$



(C. Zhu, S. R. Hudson, Y. Song, Y. Wan 2017)

# Motivation: How much does including finite coil thickness change optimized coils?



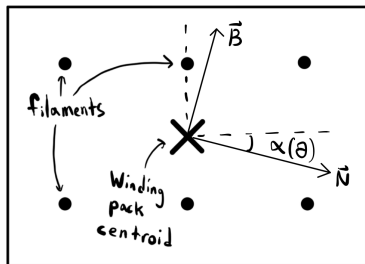
## Deviations in $B$ are second-order in the coil thickness

A simple calculation (not shown) shows that for a coil of thickness  $\delta$  a distance  $L$  from the plasma, the change in magnetic field  $\delta B$  is of order  $\delta^2/L^2$ .

$$\delta B \sim \frac{\delta^2}{L^2}$$

# One option: multi-filament approximation to winding pack

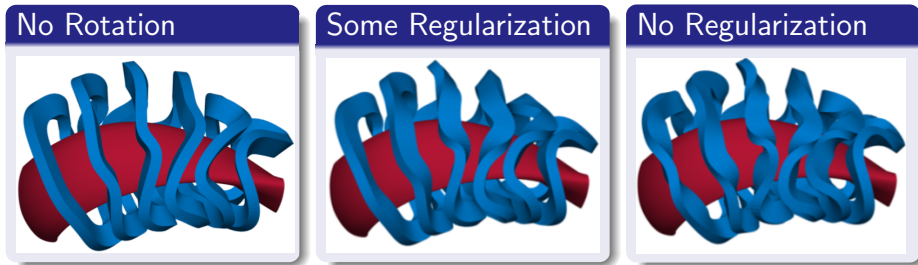
- Define frame around coil centroid
- Rotate frame by angle  $\alpha(\theta)$
- Place filaments on grid defined by rotated frame



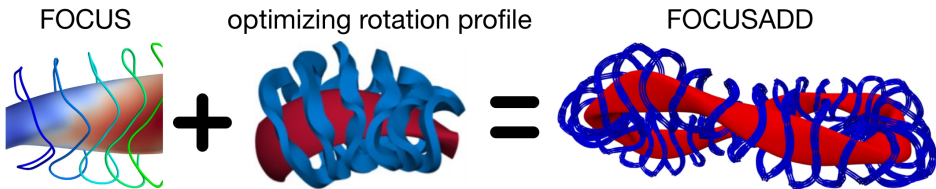
$$\alpha^i(\theta) = \sum_{m=0}^{N_{FR}-1} A_{cm}^i \cos(m\theta) + A_{sm}^i \sin(m\theta)$$

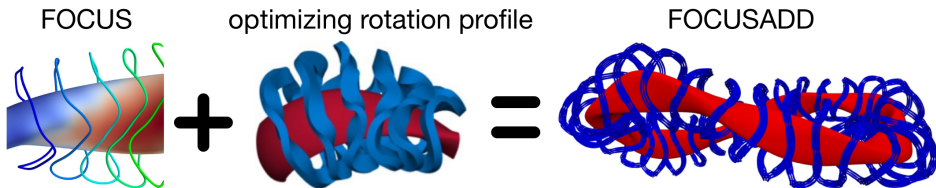
# Optimizing rotation profile of finite-build coils: HSX

Work done by Wisconsin group



L. Singh, T. Kruger et al., *Optimization of Finite-build Stellarator Coils* (Accepted to JPP)



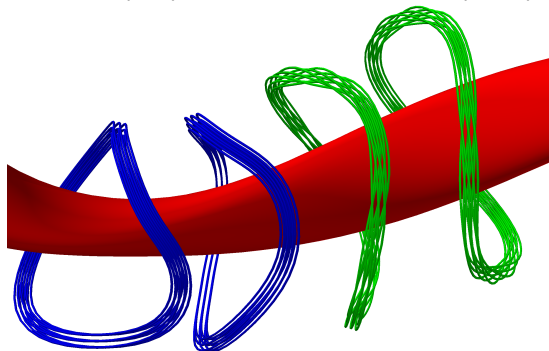


## FOCUSADD efficiently optimizes two quantities

- FOCUS optimizes the filamentary coil **centroid positions**  $x^i(\theta)$ ,  $y^i(\theta)$ ,  $z^i(\theta)$
- The Wisconsin group optimizes the **rotation profile**  $\alpha^i(\theta)$
- FOCUSADD optimizes **both** quantities at once w/ automatic differentiation

# Optimized rotation profiles with FOCUSADD

Optimized finite-build coils for rotating elliptical stellarator, with **no rotation** (left) and **no regularization** (right).



*Would anyone ever build **this**?*

Probably not. Likely the rotation profile will be set by engineering constraints.



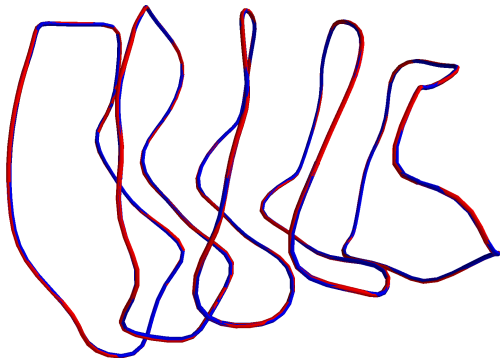
# Finite-build: I'm mostly interested in ...

- 1  $\|\Delta r\|_1$ , i.e. how much do the coils shift?
- 2  $\Delta f$ , i.e. how much does the quadratic flux change?

$\|\Delta r\|_1 \equiv$  the L1 norm of the difference between the positions of the centroid of the filamentary optimized coils and the centroid of the finite-build optimized coils.

$\Delta f \equiv$  the difference between the quadratic flux of finite-build coils which are optimized using a filamentary approximation, and the quadratic flux of finite-build coils which are optimized using the finite-build model.

# Does finite-build matter? (W7-X)



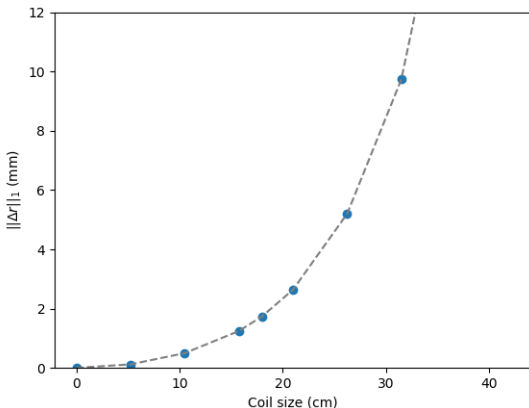
Coil centroids for optimized **filamentary** and **finite-build** (no rotation) W7-X coils.

$$\|\Delta r\|_1 = 2.6\text{mm}$$

$$\Delta f = \frac{1.45 - .45}{.45} = 240\%$$

# $\|\Delta r\|_1$ : how much do coil positions shift with finite build?

W7-X

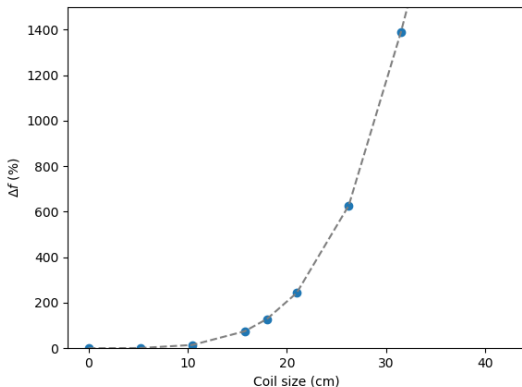


$\|\Delta r\|_1 \equiv$  the L1 norm of the difference between the positions of the centroid of the filamentary optimized coils and the centroid of the finite-build optimized coils.

W7-X coil dimensions: 17.8 by 21.2cm

# $\Delta f$ : what is the effect on the quadratic flux?

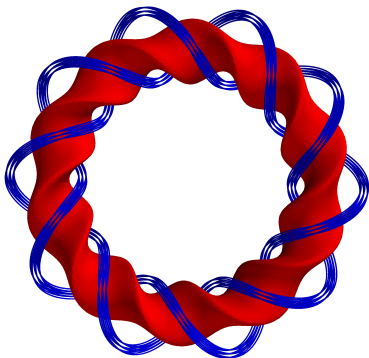
W7-X



$\Delta f \equiv$  the difference between the quadratic flux of finite-build coils which are optimized using a filamentary approximation, and the quadratic flux of finite-build coils which are optimized for their finite-build.

W7-X coil dimensions: 17.8 by 21.2cm

# Does finite-build matter? (LHD)

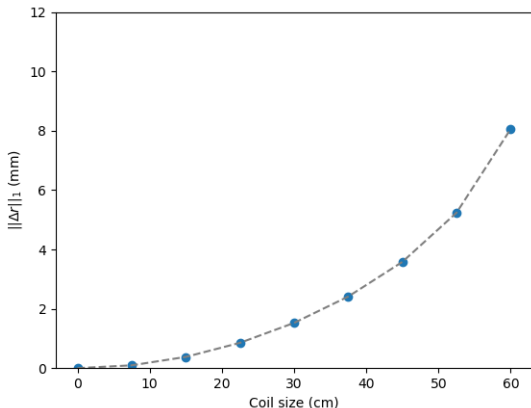


LHD finite-build coils: 30cm by 30cm. Real coils are 25cm by 45cm.

$$\|\Delta r\|_1 = 1.5\text{mm}$$

$$\Delta f = \frac{131 - 125}{125} = 5.0\%$$

# $\|\Delta r\|_1$ : how much do coil positions shift with finite build? LHD

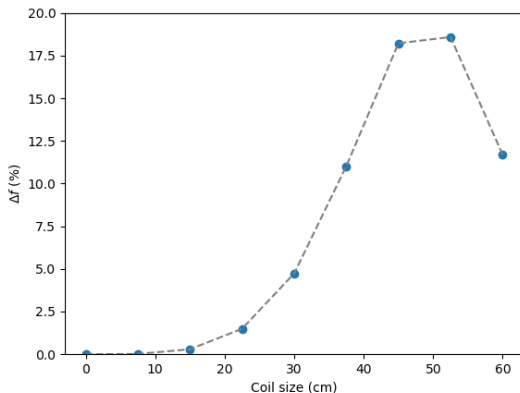


$\|\Delta r\|_1 \equiv$  the L1 norm of the difference between the positions of the centroid of the filamentary optimized coils and the centroid of the finite-build optimized coils.

LHD coil dimensions: 45 by 25cm

# $\Delta f$ : what is the effect on the quadratic flux?

LHD



$\Delta f \equiv$  the difference between the quadratic flux of finite-build coils which are optimized using a filamentary approximation, and the quadratic flux of finite-build coils which are optimized for their finite-build.

LHD coil dimensions: 45 by 25cm

# In conclusion, finite-build can be important when

- Coil tolerances are small
- Coil-plasma distance is small
- Coils are large



# Part II: Discussion of automatic differentiation (AD)

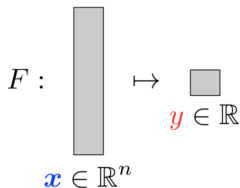


*Stan Math Library*

*differentiable C++ for linear algebra & probability*

# Taking a derivative

$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$

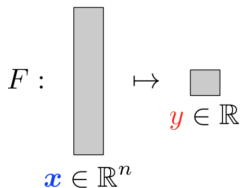


$$F = D \circ C \circ B \circ A \qquad y = F(\mathbf{x}) = D(C(B(A(\mathbf{x}))))$$

$$y = D(\mathbf{c}), \quad \mathbf{c} = C(\mathbf{b}), \quad \mathbf{b} = B(\mathbf{a}), \quad \mathbf{a} = A(\mathbf{x})$$

# Taking a derivative

$$F : \mathbb{R}^n \rightarrow \mathbb{R}$$



$$F = D \circ C \circ B \circ A \qquad y = F(\mathbf{x}) = D(C(B(A(\mathbf{x}))))$$

$$y = D(\mathbf{c}), \quad \mathbf{c} = C(\mathbf{b}), \quad \mathbf{b} = B(\mathbf{a}), \quad \mathbf{a} = A(\mathbf{x})$$

## Primitive operations

The functions  $A$ ,  $B$ ,  $C$ , and  $D$  should be understood as “primitive operations”. These could be simple primitives like `sin`, `exp`, `div`, etc, but they could also be complicated primitives like `fft`, `odeint`, `solve` or even custom primitives like `nicksfunc`.

$$y = D(\mathbf{c}), \quad \mathbf{c} = C(\mathbf{b}), \quad \mathbf{b} = B(\mathbf{a}), \quad \mathbf{a} = A(\mathbf{x})$$

$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{x}} = \left[ \frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

$$F'(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

$$\frac{\partial y}{\partial \mathbf{c}} = D'(\mathbf{c})$$

$$\frac{\partial \mathbf{c}}{\partial \mathbf{b}} = C'(\mathbf{b})$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{a}} = B'(\mathbf{a})$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} = A'(\mathbf{x})$$



$$y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

$$F'(x) = \frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

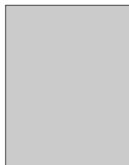
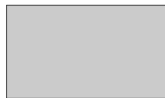
$$F'(x) = \frac{\partial y}{\partial c} \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial x}$$

$$\frac{\partial y}{\partial c} = D'(c)$$

$$\frac{\partial c}{\partial b} = C'(b)$$

$$\frac{\partial b}{\partial a} = B'(a)$$

$$\frac{\partial a}{\partial x} = A'(x)$$



## Order matters!

If  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$  takes time  $\mathcal{O}(T)$  to compute, then multiplying from right to left takes time  $\mathcal{O}(nT)$  and multiplying from left to right takes time  $\mathcal{O}(mT)$ .

# Jacobian-vector product (JVP)

A Jacobian-vector (JVP) product with one column of an identity matrix gives one column of the Jacobian matrix. E.g., for a function  $\mathbb{R}^n \rightarrow \mathbb{R}^m$ , this is

$$\mathbf{F}'(\mathbf{x})\mathbf{v} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \\ \vdots \\ \frac{\partial y_m}{\partial x_1} \end{bmatrix}$$

E.g., for a scalar function  $\mathbb{R}^n \rightarrow \mathbb{R}$ , a JVP gives a scalar.

$$\mathbf{F}'(\mathbf{x})\mathbf{v} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \cdots & \frac{\partial y}{\partial x_n} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \frac{\partial y}{\partial x_1}$$

# Vector-Jacobian product (VJP)

A vector-Jacobian product (VJP) with one row of an identity matrix gives one row of the Jacobian matrix. E.g., for a function  $\mathbb{R}^n \rightarrow \mathbb{R}^m$ , this is

$$\mathbf{v}^T \mathbf{F}'(\mathbf{x}) = [1 \ 0 \ \dots \ 0] \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \\ \vdots \\ \frac{\partial y_1}{\partial x_n} \end{bmatrix}$$

E.g., for a scalar function  $\mathbb{R}^n \rightarrow \mathbb{R}$ , a VJP gives the gradient.

$$\mathbf{v}^T \mathbf{F}'(\mathbf{x}) = [1] \begin{bmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

# AD: Forward and Reverse Modes

Forward mode: Jacobian-vector products (JVPs), build Jacobian one column at a time.  $\mathbf{v}$  is a “tangent vector”.

$$F'(\mathbf{x}) \mathbf{v} = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \left( \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \mathbf{v} \right) \right) \right)$$

Reverse mode: vector-Jacobian products (VJPs), build Jacobian one row at a time.  $\mathbf{v}$  is a “cotangent vector”.

$$\mathbf{v}^\top F'(\mathbf{x}) = \left( \left( \left( \left( \mathbf{v}^\top \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \right) \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right)$$

These are the “pushforward” and “pullback” maps from differential geometry, composed with tangent and cotangent vectors respectively.



# AD: Forward and Reverse Modes

Forward mode: Jacobian-vector products (JVPs), build Jacobian one column at a time.  $\mathbf{v}$  is a “tangent vector”.

$$F'(\mathbf{x}) \mathbf{v} = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \left( \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \mathbf{v} \right) \right) \right)$$

Reverse mode: vector-Jacobian products (VJPs), build Jacobian one row at a time.  $\mathbf{v}$  is a “cotangent vector”.

$$\mathbf{v}^\top F'(\mathbf{x}) = \left( \left( \left( \left( \mathbf{v}^\top \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \right) \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \right)$$

These are the “pushforward” and “pullback” maps from differential geometry, composed with tangent and cotangent vectors respectively.

## What makes it “automatic”?

The key is that your AD library knows how to compute Jacobian-products for each primitive operation. This allows it to compose JVPs and VJPs for each primitive and compute the total JVP/VJP.

# AD: Forward and Reverse Modes

Forward mode: Jacobian-vector products (JVPs), build Jacobian one column at a time.  $\mathbf{v}$  is a “tangent vector”.

$$F'(\mathbf{x}) \mathbf{v} = \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \left( \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \left( \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \left( \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \mathbf{v} \right) \right) \right)$$

Reverse mode: vector-Jacobian products (VJPs), build Jacobian one row at a time.  $\mathbf{v}$  is a “cotangent vector”.

$$\mathbf{v}^\top F'(\mathbf{x}) = \left( \left( \left( \mathbf{v}^\top \frac{\partial \mathbf{y}}{\partial \mathbf{c}} \right) \frac{\partial \mathbf{c}}{\partial \mathbf{b}} \right) \frac{\partial \mathbf{b}}{\partial \mathbf{a}} \right) \frac{\partial \mathbf{a}}{\partial \mathbf{x}}$$

These are the “pushforward” and “pullback” maps from differential geometry, composed with tangent and cotangent vectors respectively.

## What if I want the full Jacobian, not a JVP or VJP?

AD tools can compute the full Jacobian as well, as well as higher derivatives. However, JVPs and VJPs are the building blocks of AD tools. I highly recommend the “JAX autodiff cookbook” for further explanation.

# Fundamentals of AD: A summary

- Suppose  $\mathbf{y} = f(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ .
- If  $f$  is composed of primitive operations implemented by the AD tool, then if I give it an  $\mathbf{x}$ , AD computes the numerical value of the derivative at  $\mathbf{x}$ .
- Forward mode AD, by computing JVPs, can compute the Jacobian of  $f$  in time  $\mathcal{O}(nT)$ .
- Reverse mode AD, by computing VJPs, can compute the Jacobian of  $f$  in time  $\mathcal{O}(mT)$ .

# Why is AD useful?

## 1. Simplicity

- Finding and programming analytic derivatives is time-consuming and error-prone.

# Why is AD useful?

## 1. Simplicity

- Finding and programming analytic derivatives is time-consuming and error-prone.

## 2. Ideal for gradient-based optimization

- Reverse mode AD computes the gradient of a scalar function in time  $\mathcal{O}(T)$ . This is as efficient as the best analytic methods.

# Why is AD useful?

## 1. Simplicity

- Finding and programming analytic derivatives is time-consuming and error-prone.

## 2. Ideal for gradient-based optimization

- Reverse mode AD computes the gradient of a scalar function in time  $\mathcal{O}(T)$ . This is as efficient as the best analytic methods.

## 3. Effortless gradients

- Easy to rapidly prototype new ideas and objectives.

# Why would you not compute derivatives with AD?

- You have a massive legacy code which you can't rewrite.
- You've already written your code and it efficiently computes analytic derivatives, so why change it? (Next time though, use AD...)
- You need highly specialized numerical routines which are not implemented by an AD tool. (However, you can always wrap a numerical routine with a custom `Primitive` if you really need it.)
- Memory cost of reverse mode can be large (checkpointing helps reduce this)
- "The automatic generation of optimal (in terms of robustness and efficiency) adjoint versions of large-scale simulation code is one of the great open challenges in the field of High-Performance Scientific Computing." -Uwe Naumann (2011)

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .



# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

$$\frac{d\mathbf{g}}{d\mathbf{p}} = 0 = \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} + \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

$$\frac{d\mathbf{g}}{d\mathbf{p}} = 0 = \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} + \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial \mathbf{u}}{\partial \mathbf{p}} = - \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

$$\frac{d\mathbf{g}}{d\mathbf{p}} = 0 = \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} + \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial \mathbf{u}}{\partial \mathbf{p}} = - \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} - \frac{\partial f}{\partial \mathbf{u}} \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

$$\frac{d\mathbf{g}}{d\mathbf{p}} = 0 = \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} + \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial \mathbf{u}}{\partial \mathbf{p}} = - \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} - \frac{\partial f}{\partial \mathbf{u}} \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial f}{\partial \mathbf{u}} \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} = -\boldsymbol{\lambda}^T$$

# Review: The adjoint method

Differentiating under constraints requires a linear solve

$\mathbf{p}^* = \arg \min_{\mathbf{p}} f(\mathbf{u}(\mathbf{p}), \mathbf{p})$  s.t.  $\mathbf{g}(\mathbf{u}, \mathbf{p}) = 0$  defines  $\mathbf{u}$ . Linearize around  $\mathbf{u}$ .

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \frac{\partial f}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}}$$

$$\frac{d\mathbf{g}}{d\mathbf{p}} = 0 = \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} + \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial \mathbf{u}}{\partial \mathbf{p}} = - \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$

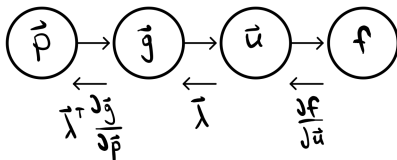
$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} - \frac{\partial f}{\partial \mathbf{u}} \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} \frac{\partial \mathbf{g}}{\partial \mathbf{p}} \Rightarrow \frac{\partial f}{\partial \mathbf{u}} \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^{-1} = -\boldsymbol{\lambda}^T$$

$$\boxed{\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \boldsymbol{\lambda}^T \frac{\partial \mathbf{g}}{\partial \mathbf{p}}} \text{ where } \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^T \boldsymbol{\lambda} = -\frac{\partial f}{\partial \mathbf{u}}$$

# Understanding the adjoint method with AD

The adjoint method computes the VJP of a constraint equation

$$\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \lambda^T \frac{\partial \mathbf{g}}{\partial \mathbf{p}}$$



- The constraint equation  $\mathbf{g} = 0$  is a primitive operation taking  $\mathbf{p}$  as input and outputting  $\mathbf{u}$ .
- $\frac{\partial f}{\partial \mathbf{u}}$  is the cotangent vector  $\mathbf{v}$  and  $\bar{\lambda}$  is the (cotangent) vector in the VJP  $\bar{\lambda} = \mathbf{v}^T \mathbf{J}$ .
- The equation  $\left(\frac{\partial \mathbf{g}}{\partial \mathbf{u}}\right)^T \bar{\lambda} = -\frac{\partial f}{\partial \mathbf{u}}$  is the adjoint equation for the primitive  $\mathbf{g} = 0$ .

# What even is an adjoint equation?

Adjoint equation  $\neq$  adjoint method

An adjoint equation is an equation which computes the VJP of a given primitive operation.

Suppose  $f_i$  is a primitive operation which takes input  $\mathbf{x}_1 \in \mathbb{R}^a$  and has output  $\mathbf{x}_2 \in \mathbb{R}^b$ . Then the adjoint equation for  $f_i$  takes a cotangent vector  $\mathbf{v}_2 \in \mathbb{R}^b$  as input and outputs a cotangent vector  $\mathbf{v}_1 = \mathbf{v}_2^T \mathbf{J}_i \in \mathbb{R}^a$ .



# What even is an adjoint equation?

Adjoint equation  $\neq$  adjoint method

An adjoint equation is an equation which computes the VJP of a given primitive operation.

Suppose  $f_i$  is a primitive operation which takes input  $\mathbf{x}_1 \in \mathbb{R}^a$  and has output  $\mathbf{x}_2 \in \mathbb{R}^b$ . Then the adjoint equation for  $f_i$  takes a cotangent vector  $\mathbf{v}_2 \in \mathbb{R}^b$  as input and outputs a cotangent vector  $\mathbf{v}_1 = \mathbf{v}_2^T \mathbf{J}_i \in \mathbb{R}^a$ .

## A different interpretation of AD via the adjoint method

Reverse mode AD libraries know how to setup and solve adjoint equations and compose their solutions for each primitive operation in the library. For example, the adjoint equation for the primitive  $y = \sin(x)$  is  $dx = dy * \cos(x)$ .



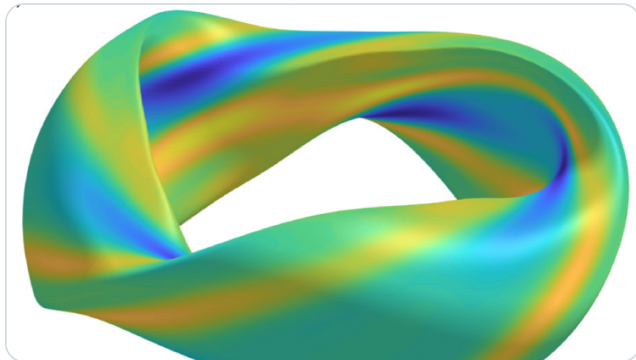
**Nick McGreivy** @NMgreivy · Jul 22



The adjoint method for PDE-constrained optimization: the conclusions of my 9-month struggle to understand the word "adjoint".

I'll be talking about how automatic differentiation (AD) can help us better understand the adjoint method, and vice versa. Let's get started!

1/



- Automatic differentiation allows the efficient computation of derivatives in stellarator coil design. This has allowed the development of a new finite-build coil design code, FOCUSADD.
- Accounting for the finite build of coils can be important when coil tolerances are small, when the coil-plasma distance is small, and/or when coil thicknesses are large.
- AD is about (i) primitive operations and (ii) JVPs/VJPs.
- AD libraries know how to setup and solve adjoint equations for each primitive operation in the library.

## Thank you

# Additional Slides

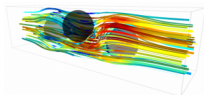
# AD of FEM code: dolfin-adjoint & FEniCS

The dolfin-adjoint project automatically derives the discrete adjoint and tangent linear models from a forward model written in the Python interface to **FEniCS** and **Firedrake**.



## dolfin-adjoint

**FEniCS** enables users to quickly translate scientific models into efficient finite element code. **Firedrake** is an automated system for the solution of partial differential equations using the finite element method (FEM).



## *Firedrake*

**dolfin-adjoint is not an AD library!**

dolfin-adjoint is similar to an AD library, but operates at a much higher level of abstraction. It works because it focuses only on FEM code.

# AD and the adjoint method for nonlinear equations

The adjoint method allows us to differentiate under the constraint  $\mathbf{g} = 0$ .

$$\boxed{\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \boldsymbol{\lambda}^T \frac{\partial \mathbf{g}}{\partial \mathbf{p}}} \text{ where } \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^T \boldsymbol{\lambda} = -\frac{\partial f}{\partial \mathbf{u}}$$

# AD and the adjoint method for nonlinear equations

The adjoint method allows us to differentiate under the constraint  $\mathbf{g} = 0$ .

$$\boxed{\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \boldsymbol{\lambda}^T \frac{\partial \mathbf{g}}{\partial \mathbf{p}}} \text{ where } \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^T \boldsymbol{\lambda} = -\frac{\partial f}{\partial \mathbf{u}}$$

Suppose the equation  $\mathbf{g} = 0$  is solved iteratively with Newton's method:

$$\mathbf{g}(\mathbf{u}^i, \mathbf{p}) + \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \cdot (\mathbf{u}^{i+1} - \mathbf{u}^i) = 0 \Rightarrow \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \delta \mathbf{u} = -\mathbf{g}(\mathbf{u}^i, \mathbf{p})$$

# AD and the adjoint method for nonlinear equations

The adjoint method allows us to differentiate under the constraint  $\mathbf{g} = 0$ .

$$\boxed{\frac{df}{d\mathbf{p}} = \frac{\partial f}{\partial \mathbf{p}} + \boldsymbol{\lambda}^T \frac{\partial \mathbf{g}}{\partial \mathbf{p}}} \text{ where } \left( \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right)^T \boldsymbol{\lambda} = -\frac{\partial f}{\partial \mathbf{u}}$$

Suppose the equation  $\mathbf{g} = 0$  is solved iteratively with Newton's method:

$$\mathbf{g}(\mathbf{u}^i, \mathbf{p}) + \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \cdot (\mathbf{u}^{i+1} - \mathbf{u}^i) = 0 \Rightarrow \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \delta \mathbf{u} = -\mathbf{g}(\mathbf{u}^i, \mathbf{p})$$

The Newton solver already computes  $\frac{\partial \mathbf{g}}{\partial \mathbf{u}}$ ! So by solving  $\mathbf{g} = 0$ , we have everything we need to perform the adjoint method. A good AD tool will set up and solve this adjoint equation automatically.



# Properties of a great AD library

- Feels like native programming
- Intuitive API
- Full set of primitive operations implemented
- Efficient linear algebra (modern implementations wrap Eigen, numpy, or JIT-compile)
- Control flow support (loops, if statements, recursion)
- Forward and reverse
- Higher-order derivatives
- GPU support
- Checkpointing
- Differentiation through linear, non-linear solves with the adjoint method
- User-defined primitives
- MPI parallelization

# JAX: Just After eXecution (Python)

JAX is Numpy and Scipy with composable function transformations:  
JIT-compile (to CPU or GPU) with `jit`, vectorize functions with `vmap`,  
SPMD parallelization with `pmap`, and automatic differentiation with  
`grad`, `jacfwd`, and `jacrev`.



# JAX: Just After eXecution (Python)

JAX is Numpy and Scipy with composable function transformations:  
JIT-compile (to CPU or GPU) with `jit`, vectorize functions with `vmap`,  
SPMD parallelization with `pmap`, and automatic differentiation with  
`grad`, `jacfwd`, and `jacrev`.



## Not just an AD library!

JAX is super useful even if you aren't doing AD! You get the simplicity of Numpy and Scipy with the speed of JIT-compilation. Programming in JAX feels just like programming in Numpy.

The Stan Math Library is a C++ template library for automatic differentiation of any order using forward, reverse, and mixed modes. It includes a range of built-in functions for probabilistic modeling, linear algebra, and equation solving.



## *Stan Math Library*

*differentiable C++ for linear algebra & probability*

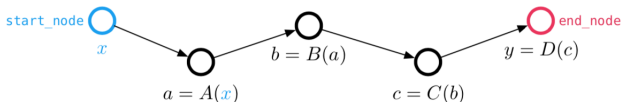
## The problem with reverse mode AD

In order to compute a vector-Jacobian product (VJP) backwards, the data to calculate each primitive VJP must be stored in memory. Storing the data from every single primitive operation cause very large memory requirements for large computations.

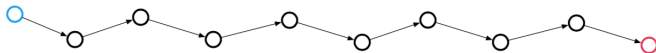
Checkpointing is a method to reduce memory requirements in exchange for increased runtime. It works by storing “checkpoints” at various points in the program and recomputing the data between checkpoints that would otherwise be stored in memory.

# Checkpointing

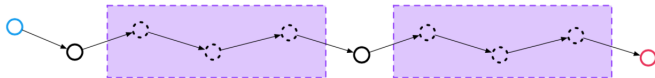
Our example function can be visualized as a graph:



Suppose we have a function which looks like this:

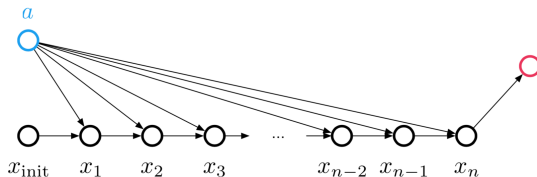


Let's use checkpointing to reduce the memory from the part of the computation in the purple boxes. We place checkpoints to the left of the purple boxes and recompute the functions in the purple boxes on the backwards pass.



# AD and the adjoint method for iterative algorithms

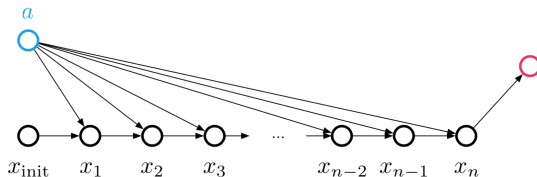
Suppose we have an iterative algorithm for solving some set of equations, and we want to compute the derivative of the solution with respect to some parameter  $a$ . Let  $x_{\text{init}}$  be an initial guess, and suppose the iterative algorithm runs  $n$  times before converging to a solution. This computation can be visualized with the following graph:



We could compute the derivative  $dx/da$  by using automatic differentiation on the entire computation. This might be very inefficient, both in terms of runtime and memory.

# AD and the adjoint method for iterative algorithms

Suppose we have an iterative algorithm for solving some set of equations, and we want to compute the derivative of the solution with respect to some parameter  $a$ . Let  $x_{\text{init}}$  be an initial guess, and suppose the iterative algorithm runs  $n$  times before converging to a solution. This computation can be visualized with the following graph:



We could compute the derivative  $dx/da$  by using automatic differentiation on the entire computation. This might be very inefficient, both in terms of runtime and memory.

## A clever trick

We can use the mathematical structure of the iteration to more efficiently compute the derivative. The key is that the derivative doesn't depend on  $x_{\text{init}}$ . This means that when computing the derivative, we can simply rerun the iteration with  $x_{\text{init}} = x_n$ , and compute the derivative of a single iteration of the algorithm. This is the trick for applying the adjoint method to systems of non-linear equations.

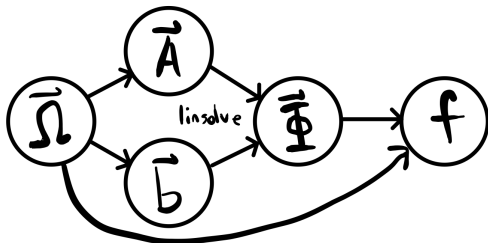


# Review: AD and the adjoint method for linear equations

Goal:  $\Omega^* = \arg \min_{\Omega} f(\Phi(\Omega), \Omega)$

Use GD:  $\Omega^{n+1} = \Omega^n - \eta \frac{\partial f}{\partial \Omega}$

s.t.  $\mathbf{A}(\Omega)\Phi = \mathbf{b}(\Omega)$

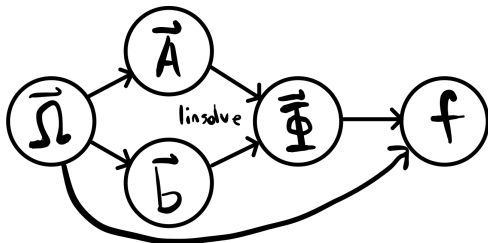


# Review: AD and the adjoint method for linear equations

Goal:  $\Omega^* = \arg \min_{\Omega} f(\Phi(\Omega), \Omega)$

Use GD:  $\Omega^{n+1} = \Omega^n - \eta \frac{\partial f}{\partial \Omega}$

s.t.  $\mathbf{A}(\Omega)\Phi = \mathbf{b}(\Omega)$



Using the adjoint method, we have

$$\boxed{\frac{df}{d\Omega} = \frac{\partial f}{\partial \Omega} + \lambda^T \frac{\partial \mathbf{b}}{\partial \Omega} - \lambda^T \frac{\partial \mathbf{A}}{\partial \Omega} \Phi} \quad \text{where } \mathbf{A}^T \lambda = \frac{\partial f}{\partial \Phi}$$

AD tools setup and solve the adjoint equation of a linear system automatically, e.g. JAX uses `np.linalg.solve` and `grad`.

# Should AD be partially adopted by the stellarator optimization community? It's worth considering.

This is an enormously complex question. I'd love to discuss this in depth after the talk or offline.

# Should AD be partially adopted by the stellarator optimization community? It's worth considering.

This is an enormously complex question. I'd love to discuss this in depth after the talk or offline.

## Pros

- Gradient-based optimization of high-dimensional non-convex objective functions has been successful in many domains.
- AD and the adjoint method work together particularly well.
- If  $N = 50$ , does a factor of 10-20 increase in computational speed matter?
- Exact derivative needed?
- Much easier to rewrite an existing code we understand than write a new code.

# Should AD be partially adopted by the stellarator optimization community? It's worth considering.

This is an enormously complex question. I'd love to discuss this in depth after the talk or offline.

## Pros

- Gradient-based optimization of high-dimensional non-convex objective functions has been successful in many domains.
- AD and the adjoint method work together particularly well.
- If  $N = 50$ , does a factor of 10-20 increase in computational speed matter?
- Exact derivative needed?
- Much easier to rewrite an existing code we understand than write a new code.

## Cons

- Is there sufficient demand for rewriting STELLOPT with an AD tool? Does our team have the right expertise?
- The stellarator community seems to like FORTRAN. Bad for AD.
- Does the right tool exist?
- Is gradient-free Bayesian optimization better? Bayesian optimization with gradients? Do we have resources to try all the above?
- Memory manageable?

# AD in machine learning (ML):

A 2009 blog post made a convincing argument that ML researchers should use AD.

## Criticisms of blog post:

- Computing derivatives distracts from what the field actually wants to accomplish.
- Valuable researcher time is wasted.
- Leads to preference for functions they are capable of manually deriving gradients for.

AD was eventually fully adopted by the ML community.

# Structure of FOCUSADD (1 of 6)

Coil centroid is parametrized in free space with a Fourier series, as in FOCUS. Here  $\mathbf{r}_i(\theta)$  is the position of the  $i$  coil centroid.

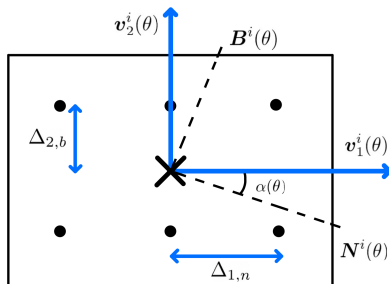
$$x^i(\theta) = \sum_{m=0}^{N_F-1} X_{cm}^i \cos(m\theta) + X_{sm}^i \sin(m\theta)$$

$$y^i(\theta) = \sum_{m=0}^{N_F-1} Y_{cm}^i \cos(m\theta) + Y_{sm}^i \sin(m\theta)$$

$$z^i(\theta) = \sum_{m=0}^{N_F-1} Z_{cm}^i \cos(m\theta) + Z_{sm}^i \sin(m\theta)$$

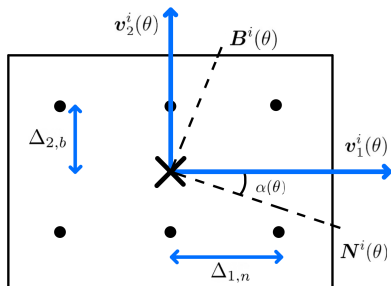
## Structure of FOCUSADD (2 of 6)

The multi-filament winding pack surrounds the coil centroid. For the  $i$ th coil, the axes of the winding pack  $\mathbf{v}_1^i$  and  $\mathbf{v}_2^i$  are rotated by an angle  $\alpha^i$  relative to the normal  $\mathbf{N}^i$  and binormal  $\mathbf{B}^i$  vectors. The normal vector is defined as the component perpendicular to the tangent of the vector from the coil center-of-mass to the local point.





## Structure of FOCUSADD (3 of 6)



$\alpha$  is parametrized by another Fourier series, giving the coil the freedom to twist in space.

$$\alpha^i(\theta) = \frac{N_R \theta}{2} + \sum_{m=0}^{N_{FR}-1} A_{cm}^i \cos(m\theta) + A_{sm}^i \sin(m\theta)$$

## Structure of FOCUSADD (4 of 6)

Once we have  $\mathbf{v}_1^i$  and  $\mathbf{v}_2^i$ , we can compute the position of the  $N_n \times N_b$  filaments for each coil. We do this using the following formula for the  $n$ th and  $b$ th filaments, where  $n$  runs from 0 to  $N_n - 1$  and  $b$  runs from 0 to  $N_b - 1$ . Here,  $l_n$  is the spacing between the filaments in the  $\mathbf{v}_1$  direction, and  $l_b$  is the spacing between the filaments in the  $\mathbf{v}_2$  direction.

$$\mathbf{r}_{n,b}^i(\theta) = \mathbf{r}_{central}^i + \left[ n - \frac{N_n - 1}{2} \right] l_n \mathbf{v}_1^i(\theta) + \left[ b - \frac{N_b - 1}{2} \right] l_b \mathbf{v}_2^i(\theta)$$

So far we've only computed vacuum fields, using the Biot-Savart law.

$$\mathbf{B}(\mathbf{r}) = \sum_{i=1}^{N_c} \sum_{n=1}^{N_1} \sum_{b=1}^{N_2} \mu_0 I_{n,b}^i \oint \frac{d\mathbf{l}_{n,b}^i \times (\mathbf{r} - \mathbf{r}_{n,b}^i)}{|\mathbf{r} - \mathbf{r}_{n,b}^i|^3}$$

## Structure of FOCUSADD (5 of 6)

The optimization is performed using gradient descent on an objective function  $f_{total}$ , given by a sum of physics objectives and engineering objectives.

$$f_{total}(\mathbf{p}) = f_{Phys}(\mathbf{p}) + \lambda_{Eng} f_{Eng}(\mathbf{p})$$

The simplest possible physics objective was chosen, the squared surface-normal magnetic field integrated over the surface.

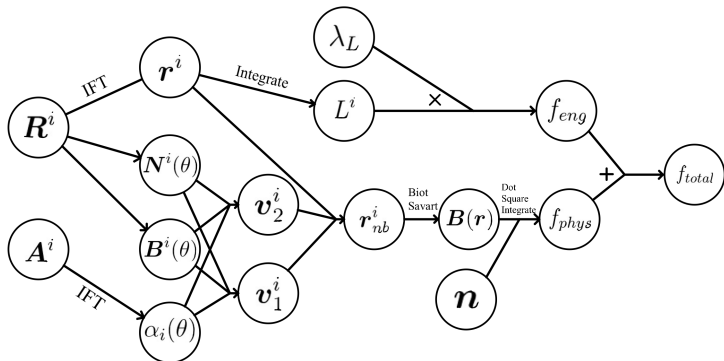
$$f_{Phys} \equiv \int_S (\mathbf{B} \cdot \mathbf{n})^2 dA$$

A simple engineering objective was chosen, namely the total length of the coils.

$$f_{Eng} \equiv \sum_{i=1}^N L_i$$

# Structure of FOCUSADD (6 of 6)

The following computational graph describes the structure of the computation performed by FOCUSADD. In my AD tool (JAX), I simply compute  $f_{total}$ , then type “grad” to get the gradient.



# Important facts about AD

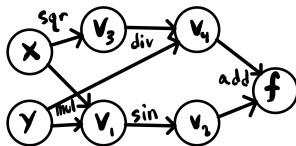
- Suppose  $\mathbf{y} = f(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ . Then first-order AD computes the numerical value of the Jacobian  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$  at a particular value of  $\mathbf{x}$ .
- Suppose  $f$  takes time  $T$  to compute.
- Forward mode AD computes the Jacobian of  $f$  in time  $\mathcal{O}(nT)$ .
- Reverse mode AD computes the Jacobian of  $f$  in time  $\mathcal{O}(mT)$ .

# How does AD work?

## Step 1: Compute function

The AD tool computes the function  $f$ , one elementary operation at a time. A representation of the 'computational graph' is built.

$$f(x, y) = \sin(xy) + x^2/y$$



## Step 2: Compute derivatives in the reverse order

Compute the derivative of the output of  $f$  with respect to each variable  $v_i$  by traversing the graph in reverse order. This is exactly the chain rule, applied in a clever way.

$$\frac{\partial f}{\partial v_i} = \sum_{\substack{j \in \text{children} \\ \text{of } i}} \frac{\partial f}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

# The details

Step 1:

$$x = 2.0$$

$$y = 3.0$$

$$v_1 = x * y = 6.0$$

$$v_2 = \sin(v_1) = -0.297$$

$$v_3 = x^2 = 4.0$$

$$v_4 = v_3 / y = 1.333$$

$$f = v_2 + v_4 = 1.054$$

Step 2:

$$\frac{\partial f}{\partial f} = 1.0$$

$$\frac{\partial f}{\partial v_4} = 1.0$$

$$\frac{\partial f}{\partial v_3} = \frac{\partial f}{\partial v_4} \frac{\partial v_4}{\partial v_3} = 0.333$$

$$\frac{\partial f}{\partial v_2} = 1.0$$

$$\frac{\partial f}{\partial v_1} = \frac{\partial f}{\partial v_2} \frac{\partial v_2}{\partial v_1} = 0.960$$

Step 2, continued:

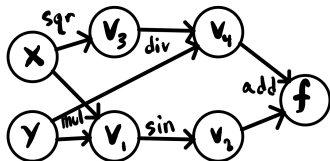
$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial v_1} \frac{\partial v_1}{\partial y} + \frac{\partial f}{\partial v_4} \frac{\partial v_4}{\partial y}$$

$$\frac{\partial f}{\partial y} = 1.476$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial v_3} \frac{\partial v_3}{\partial x} + \frac{\partial f}{\partial v_1} \frac{\partial v_1}{\partial x}$$

$$\frac{\partial f}{\partial x} = 4.214$$

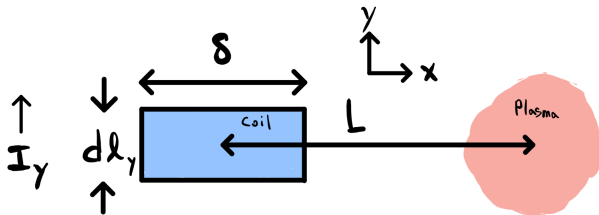
$$f(x, y) = \sin(xy) + x^2/y$$



$$\frac{\partial f}{\partial v_i} = \sum_{j \in \text{children of } i} \frac{\partial f}{\partial v_j} \frac{\partial v_j}{\partial v_i}$$

# Deviations in the magnetic field are **second-order** in the coil thickness divided by the coil-plasma distance

Suppose we have a coil carrying current in the  $y$ -direction with thickness  $\delta$  a distance  $L$  away from our plasma.



$$\text{Biot-Savart: } \frac{dB_z}{d\ell_y} = -\frac{\mu_0 I_y}{4\pi} \int_{-\delta/2}^{\delta/2} \frac{dx}{(L+x)^2}$$

$$dB_z \approx -\frac{\mu_0 I_y d\ell_y}{4\pi L^2} \int_{-\delta/2}^{\delta/2} \left(1 - \frac{2x}{L} + \frac{3x^2}{L^2}\right) dx \approx dB_{\text{filament}} \left(1 + \frac{\delta^2}{4L^2}\right)$$