# CS 4803 / 7643: Deep Learning

Topics:

– Automatic Differentiation

  – (Finish) Forward mode vs Reverse mode AD

  – Patterns in backprop

  – Jacobians in FC+ReLU NNs

Dhruv Batra

Georgia Tech

# Administrativia

- HW1 Reminder
  – Due: 09/26, 11:55pm

# Project

- Goal
  - Chance to take on something open-ended
  - Encouraged to apply to your research (computer vision, NLP, robotics,…)

- Main categories
  - Application/Survey
    - Compare a collection of existing algorithms on a new application domain of your interest
  - Formulation/Development
    - Formulate a new model or algorithm for a new or old problem
  - Theory
    - Theoretically analyze an existing algorithm

# Project

- Rules
  - **Combine with other classes / research / credits / anything**
    - You have our blanket permission
    - Get permission from other instructors; delineate different parts
  - Must be done this semester.
  - Groups of 3-4

- Expectations
  - 20% of final grade = individual effort equivalent to 1 HW
  - Expectation scales with team size
  - Most work will be done in Nov but please plan early.

# Project Ideas

- NeurIPS Reproducibility Challenge
  - https://reproducibility-challenge.github.io/neurips2019/
  - https://reproducibility-challenge.github.io/neurips2019/task/

**Reproducibility** Challenge    NeurIPS 2019   Task Description   Resources   Registration   Important Dates   Organizers



# Reproducibility Challenge @ NeurIPS 2019

The Annual Machine Learning Reproducibility Challenge

# Computing

- **Major bottleneck**
  - GPUs

- **Options**
  - Your own / group / advisor's resources

  - Google Cloud Credits
    - $50 credits to every registered student courtesy Google

  - Google Colab
    - jupyter-notebook + free GPU instance

# Administrativia

- Project Teams Google Doc
  - https://docs.google.com/spreadsheets/d/1ouD6ctaemV_3nb2MQHs7rUOAaW9DFLu8I5Zd3yOFs7E/edit?usp=sharing
  - Project Title
  - 1-3 sentence project summary TL;DR
  - Team member names

# Recap from last time

# How do we compute gradients?

- Analytic or "Manual" Differentiation

- Symbolic Differentiation

- Numerical Differentiation

- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka "backprop"

# Chain Rule: Composite Functions

$$\left[ L(x) = f(g(x)) = (f \circ g)(x) \right]$$

$$f(x) = g_\ell(g_{\ell-1} \cdots g_1(x))$$

$$L(w) = (g_\ell \circ g_{\ell-1} \cdots \circ g_1)(x)$$

$$\frac{\partial L}{\partial w}$$

# Chain Rule: <u>Scalar</u> Case

$$x \xrightarrow{g(\cdot)} z \xrightarrow{f(\cdot)} y \to a \qquad x, y, z \in \mathbb{R}^1$$
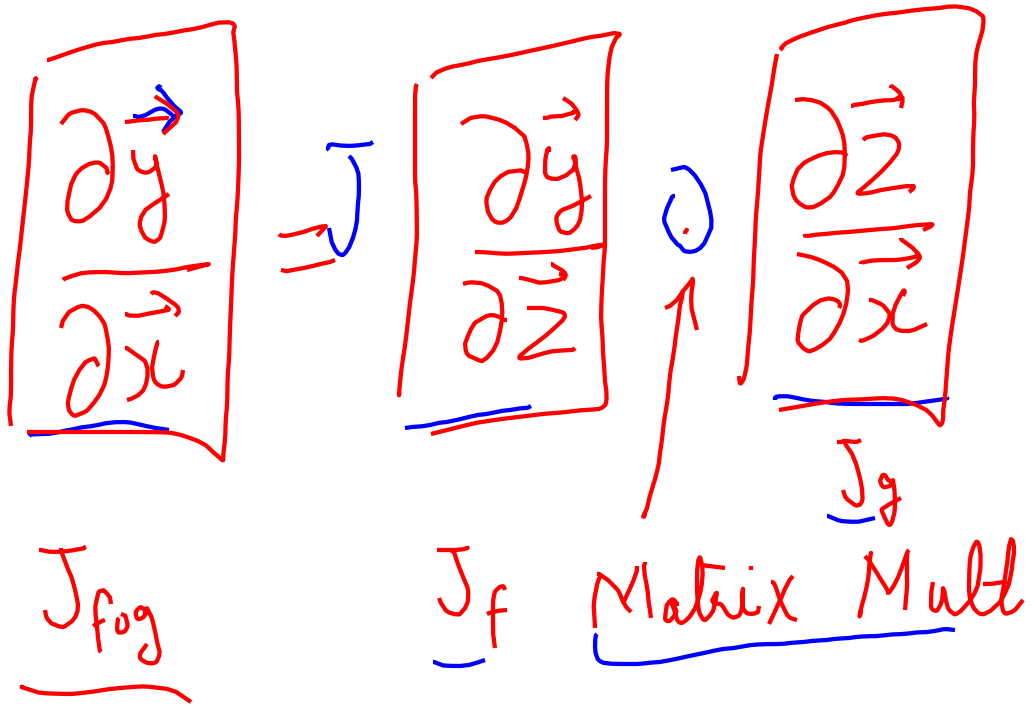
$$a \in \mathbb{R}^1$$

$$y = f(\underbrace{g(x)}_{z})$$

$$\boxed{\frac{\partial y}{\partial x}} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x}$$

Scalar prod.

$$\frac{\partial a}{\partial x} = \frac{\partial a}{\partial y} \cdot \boxed{\frac{\partial y}{\partial x}}$$

$$= \frac{\partial a}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x}$$
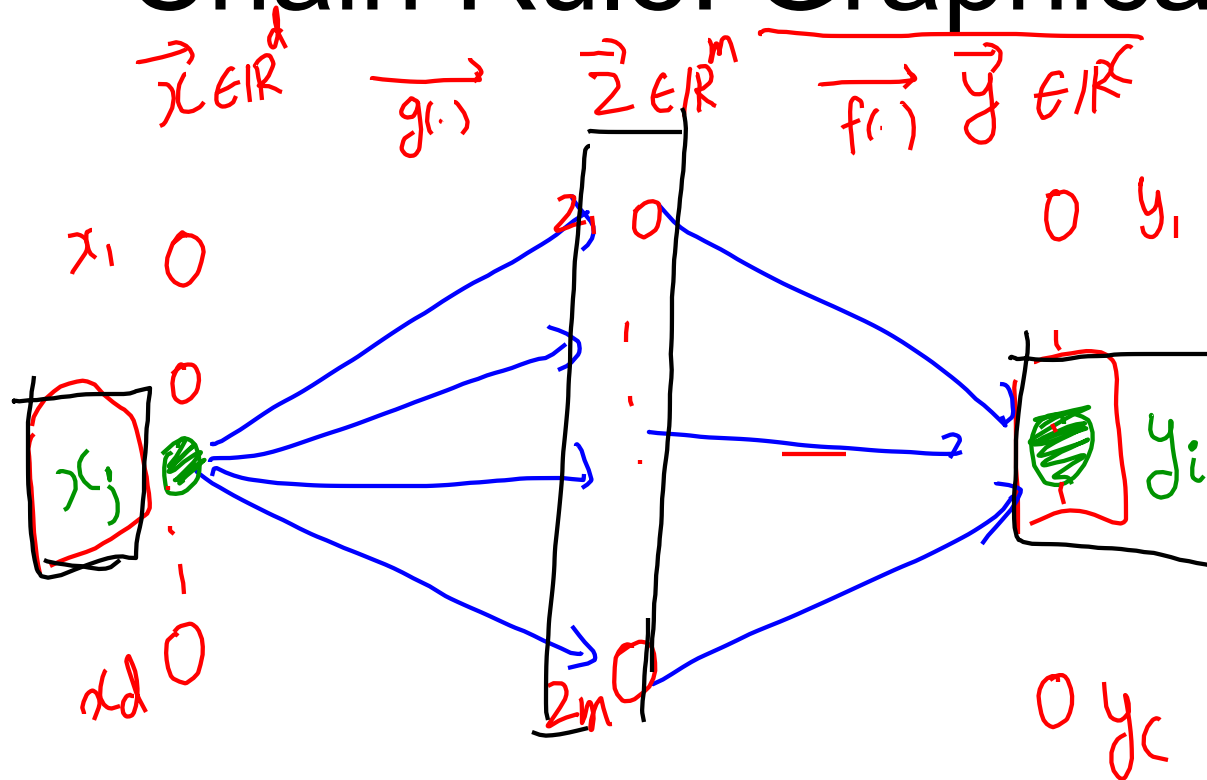
# Chain Rule: Vector Case

$$\vec{x} \in \mathbb{R}^d \xrightarrow{g(\cdot)} \vec{z} \in \mathbb{R}^m \xrightarrow{f(\cdot)} \vec{y} \in \mathbb{R}^c \rightarrow \vec{a}$$

$$\mathbb{R}^d \rightarrow \mathbb{R}^m \qquad \mathbb{R}^m \rightarrow \mathbb{R}^c$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = J = \frac{\partial \vec{y}}{\partial \vec{z}} (\cdot) \frac{\partial \vec{z}}{\partial \vec{x}}$$

$$J_{f \circ g} \qquad J_f \quad \text{Matrix Mult} \qquad J_g$$

# Chain Rule: Jacobian view

$$\frac{\partial \vec{y}}{\partial \vec{x}}$$

$$\frac{\partial \vec{y}}{\partial \vec{z}}$$

$$\frac{\partial \vec{z}}{\partial \vec{x}}$$

$$i\left[ \cdots \frac{\partial y_i}{\partial x_j} \cdots \right]_{c \times d} = i\left[ \cdots \frac{\partial y_i}{\partial z_k} \cdots \right]_{c \times m} \left[ \frac{\partial z_k}{\partial x_j} \right]_{m \times d}$$

$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$$
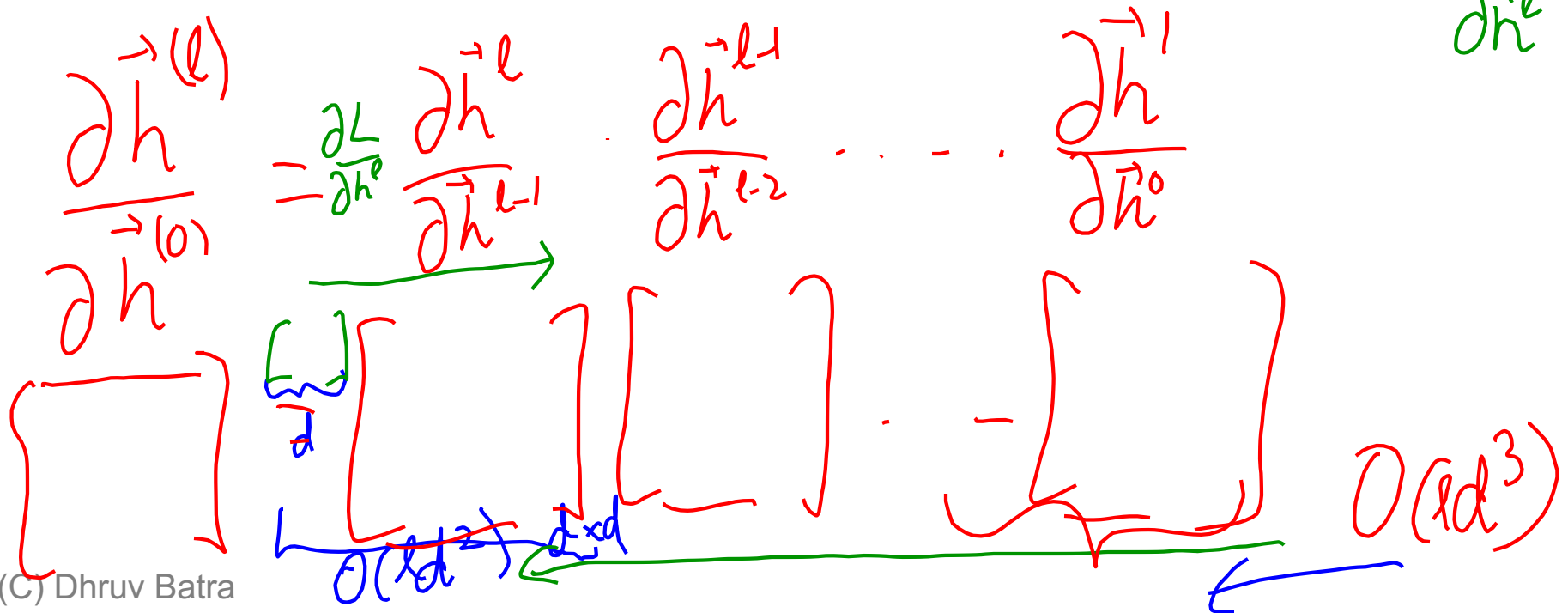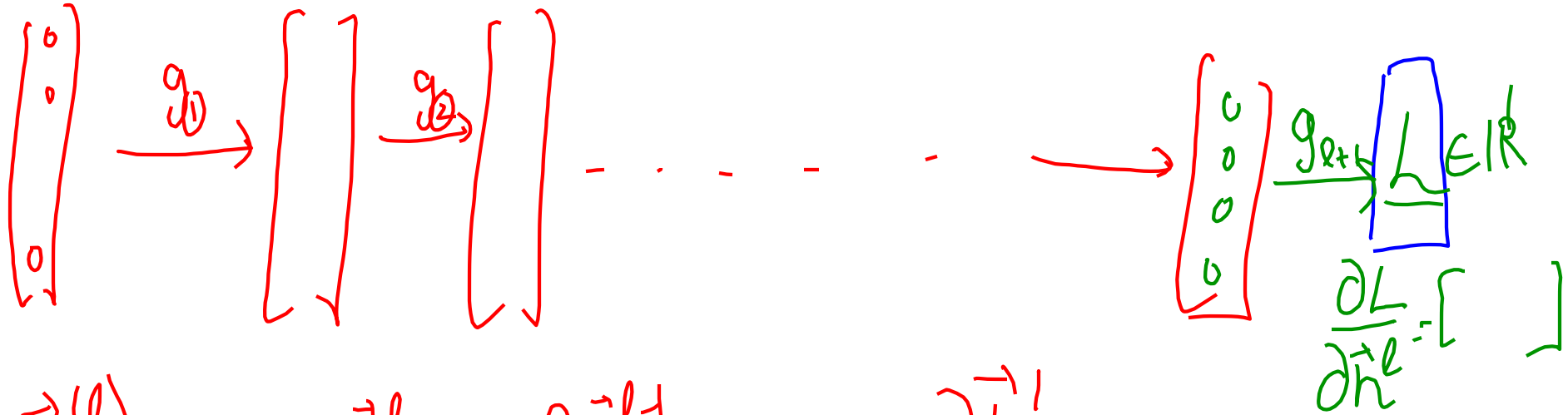
# Chain Rule: Graphical view



$$\frac{\partial y_i}{\partial x_j} = \sum_{Paths} \frac{\partial y_i}{\partial z_k} \cdot \frac{\partial z_k}{\partial x_j} \cdots$$

$k$ is on path

# Chain Rule: Cascaded

$$\vec{h}^0 \xrightarrow{g_1()} \vec{h}^{(1)} \longrightarrow \vec{h}^2 \cdots\cdots\cdots \longrightarrow \vec{h}^{(\ell)} \in \mathbb{R}^d$$

$$\vec{h}^0 \in \mathbb{R}^d$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \xrightarrow{g_1} \begin{bmatrix} \\ \\ \end{bmatrix} \xrightarrow{g_2} \begin{bmatrix} \\ \\ \end{bmatrix} \cdots\cdots\cdots \longrightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{g_{\ell+1}} L \in \mathbb{R}$$

$$\frac{\partial L}{\partial \vec{h}^\ell} = [\quad\quad]$$

$$\frac{\partial \vec{h}^{(\ell)}}{\partial \vec{h}^{(0)}} = \frac{\partial L}{\partial \vec{h}^0} \frac{\partial \vec{h}^\ell}{\partial \vec{h}^{\ell-1}} \cdot \frac{\partial \vec{h}^{\ell-1}}{\partial \vec{h}^{\ell-2}} \cdots\cdots \frac{\partial \vec{h}^1}{\partial \vec{h}^0}$$

$$\begin{bmatrix} \\ \\ \\ \end{bmatrix} \underbrace{\begin{bmatrix} \\ \\ \\ \end{bmatrix}}_{d} \begin{bmatrix} \\ \\ \\ \end{bmatrix} \cdots - \begin{bmatrix} \\ \\ \\ \end{bmatrix} \qquad O(\ell d^3)$$

$$O(\ell d^2) \quad d \times d$$

# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering

- Auto-Diff
  - A family of algorithms for

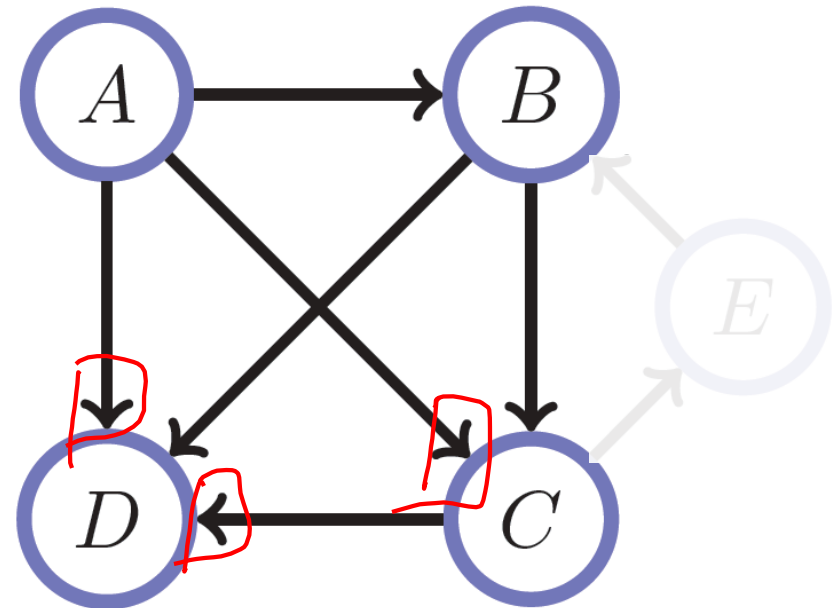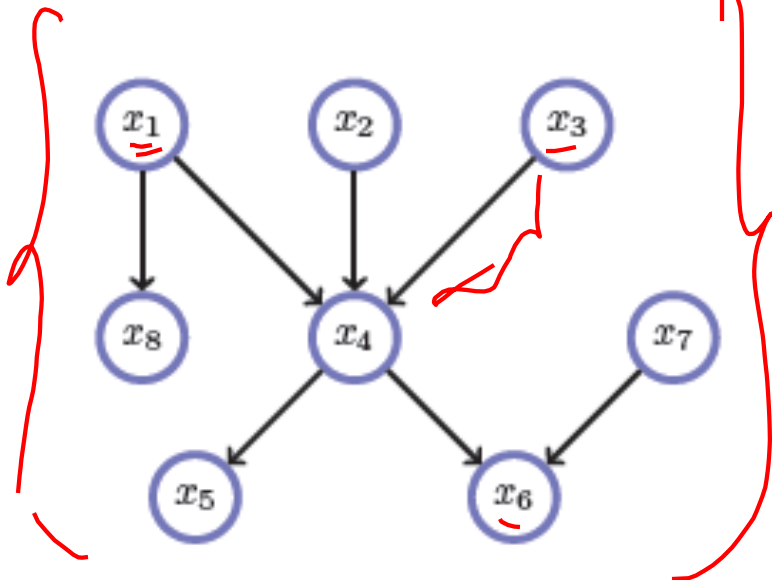    implementing chain-rule on computation graphs

# Directed Acyclic Graphs (DAGs)

- Exactly what the name suggests
  - Directed edges
  - No (directed) cycles
  - Underlying undirected cycles okay

$$G = (V, E)$$

$$E = \{(v_i, v_j) \mid v_i, v_c \in V\}$$
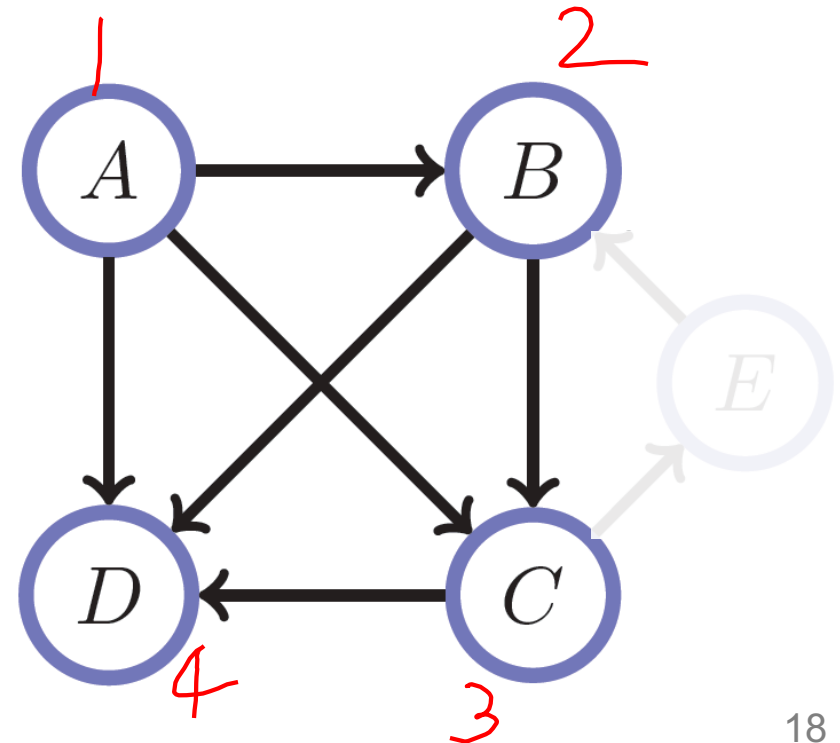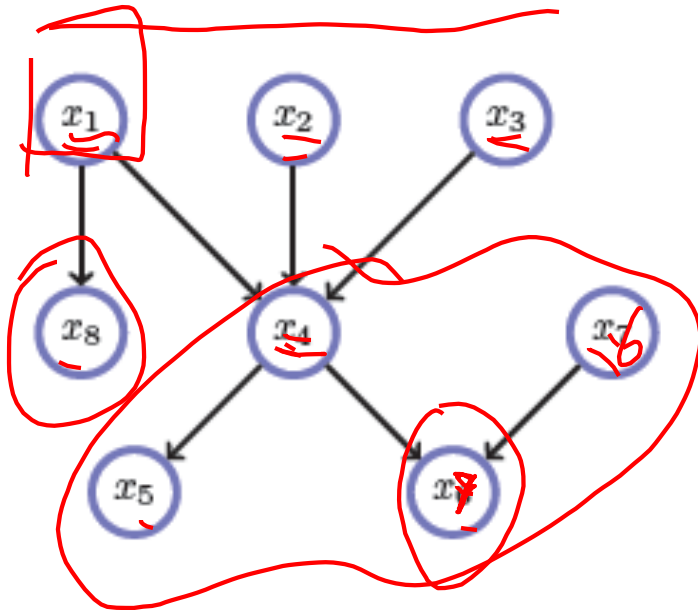
Vertex
Vertices

# Directed Acyclic Graphs (DAGs)

- Concept
  - Topological Ordering

$\exists$ bijection $\sigma : V \rightarrow \{1, \ldots, n\}$

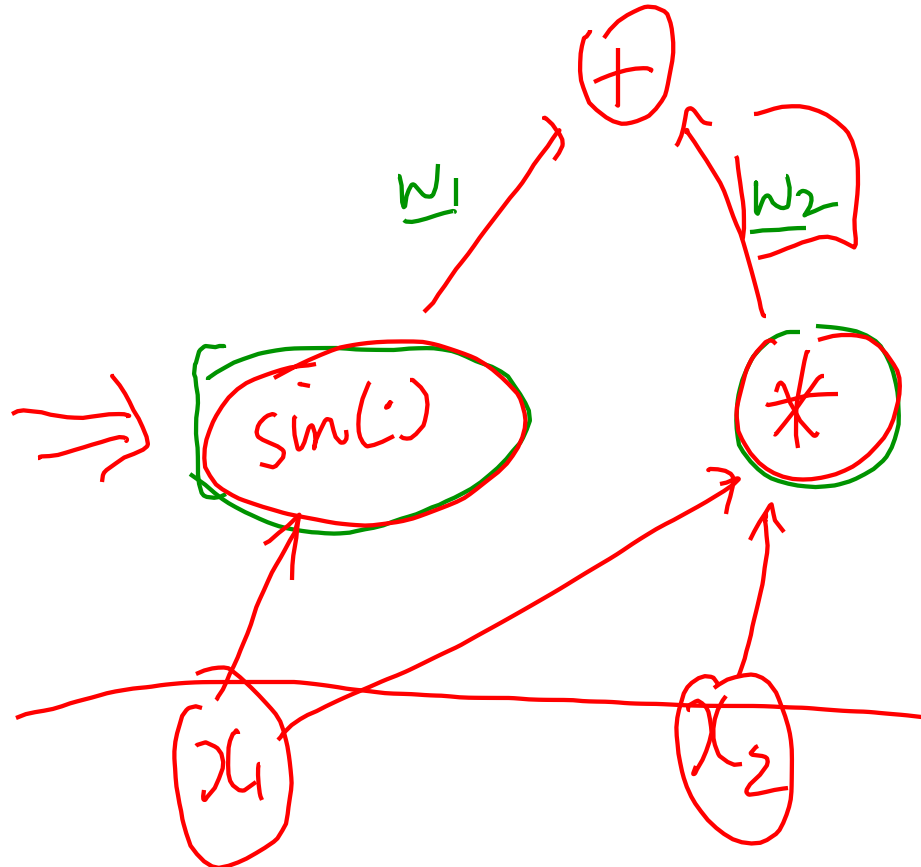s.t $\forall (v_i, v_j) \in E$

$$\sigma(v_i) < \sigma(v_j)$$

# Computational Graphs

- Notation

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

# Deep Learning = Differentiable Programming

- Computation = Graph
  - Input = Data + Parameters
  - Output = Loss
  - Scheduling = Topological ordering

- Auto-Diff
  - A family of algorithms for
    implementing chain-rule on computation graphs

# Forward mode AD

layer $l$

$\boxed{L}$

$\dfrac{\partial L}{\partial \vec{x}}$

$\vec{\dot{x}}$

$\boxed{\vec{h}^{(l-1)}}$

$g(\ )$

$\vec{h}^{(l)} = g\left(\vec{h}^{\,l-1}\right)$

Input: $\boxed{\dfrac{\partial \vec{h}^{\,l-1}}{\partial \vec{x}}}$

$\boxed{\dfrac{\partial \vec{h}^{\,l}}{\partial \vec{x}}} = \boxed{\dfrac{\partial \vec{h}^{\,l}}{\partial \vec{h}^{\,l-1}}} \cdot \boxed{\dfrac{\partial \vec{h}^{\,l-1}}{\partial \vec{x}}}$

Jacobian of $g$  Input

# Reverse mode AD

$\vec{y} = A\vec{x}$

$\dfrac{\partial \vec{y}}{\partial \vec{x}} = A$

$\vec{h}^{\ell} = W\vec{h}^{\ell-1}$

$\dfrac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}} = W$

Goal: $\dfrac{\partial L}{\partial \vec{x}}$

$\vec{x}$

$\vec{h}^{\ell-1}$

$g$

$\vec{h}^{\ell} = g(\vec{h}^{\ell-1})$

$\dfrac{\partial L}{\partial \vec{h}^{\ell}}$  Input

Output

$\dfrac{\partial L}{\partial \vec{h}^{\ell-1}} = \dfrac{\partial L}{\partial \vec{h}^{\ell}} \cdot \dfrac{\partial \vec{h}^{\ell}}{\partial \vec{h}^{\ell-1}}$

Input

Jacobian of $g$

$\dfrac{\partial L}{\partial \vec{x}}$

# Plan for Today

- Automatic Differentiation
  - (Finish) Forward mode vs Reverse mode AD
  - Backprop
  - Patterns in backprop
  - Jacobians in FC+ReLU NNs
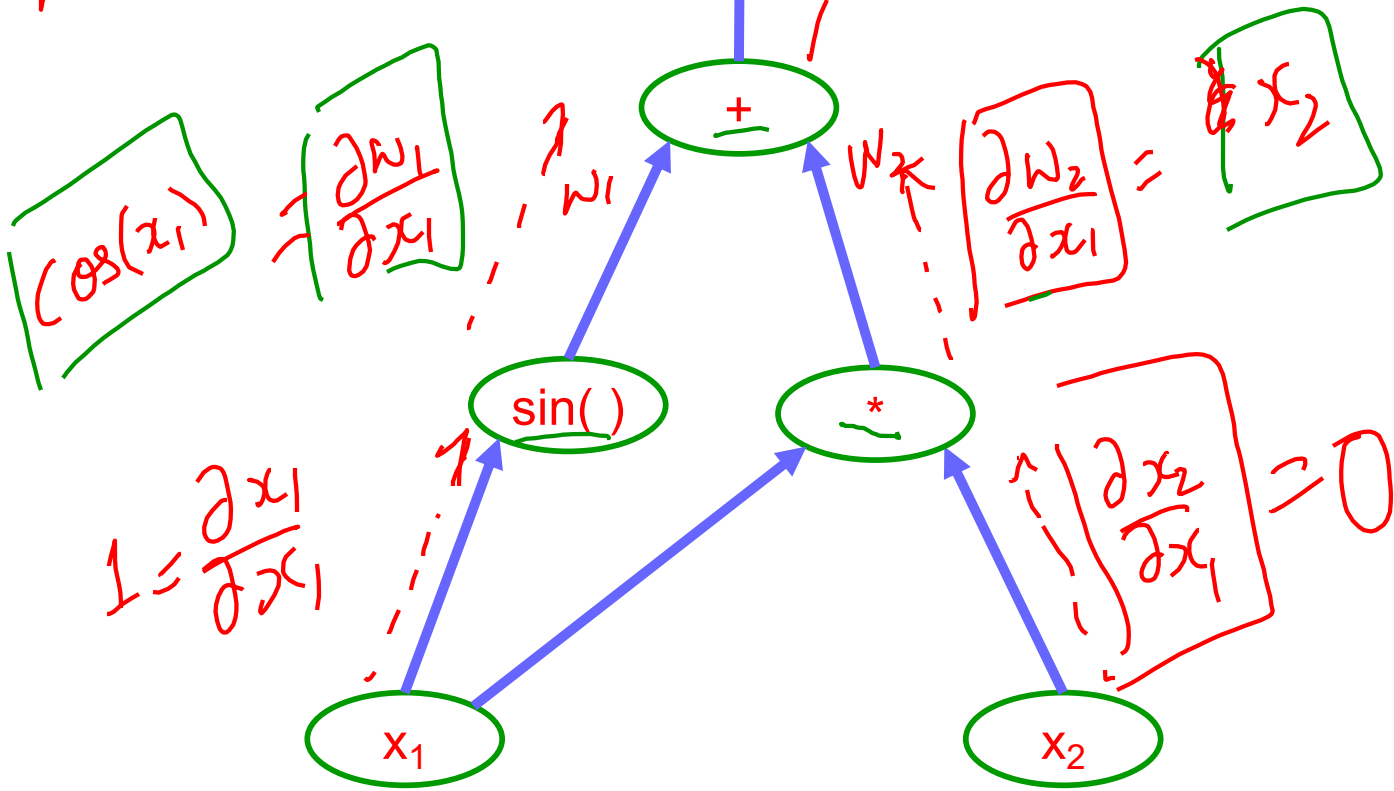
# Example: Forward mode AD

$$f(x_1, x_2) = \overbrace{x_1 x_2}^{w_2} + \overbrace{\sin(x_1)}^{w_1}$$

$$\frac{\partial f}{\partial x} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

$\frac{\partial f}{\partial x_1}$

$= \cos x_1 + x_2$

$w_1 = \sin(x_1)$

$w_1 + w_2 = w_3 \uparrow \quad \frac{\partial w_3}{\partial x_1} = \frac{\partial w_1}{\partial x_1} + \frac{\partial w_2}{\partial x_1}$
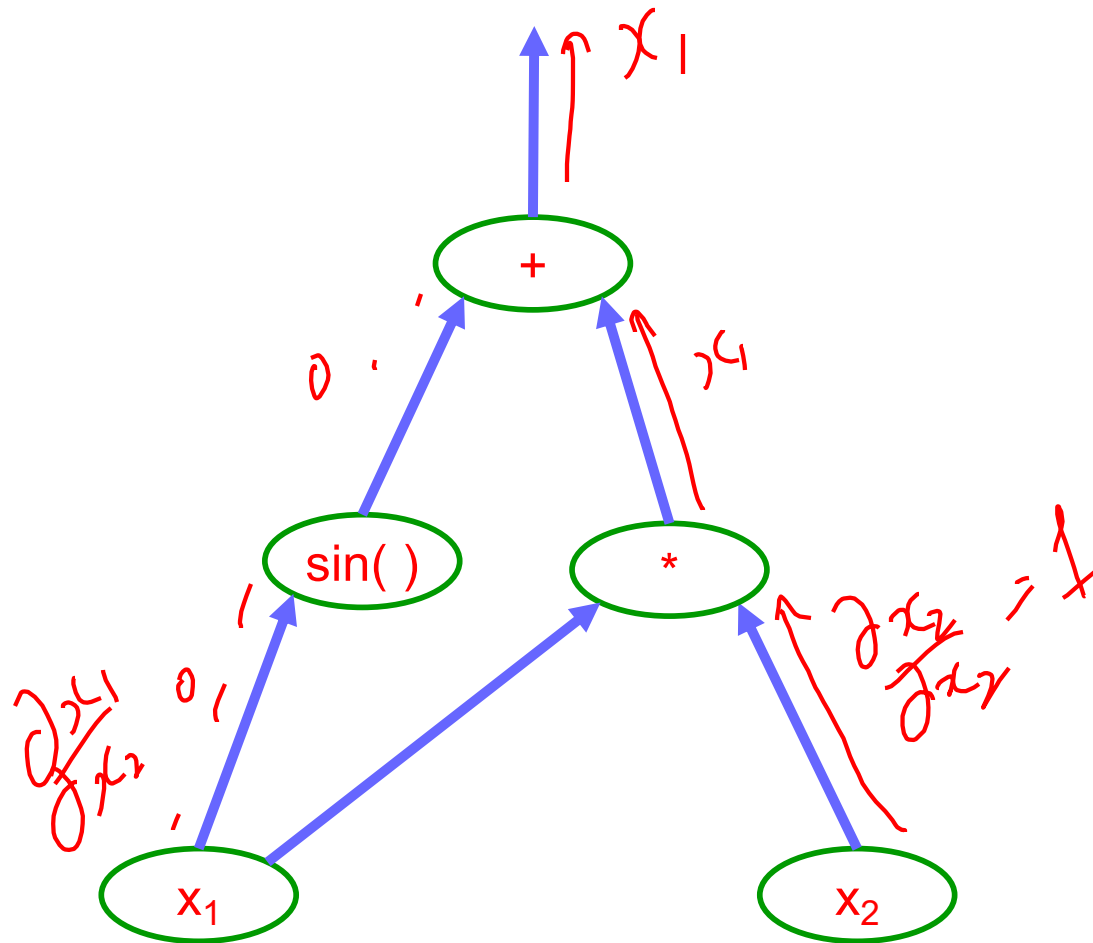
$w_2 = x_1 x_2$

$\cos(x_1) = \frac{\partial w_1}{\partial x_1}$

$w_1$

$+$

$w_2 \quad \frac{\partial w_2}{\partial x_1} = x_2$

$\sin()$

$*$

$1 = \frac{\partial x_1}{\partial x_1}$

$\frac{\partial x_2}{\partial x_1} = 0$

$x_1$

$x_2$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$
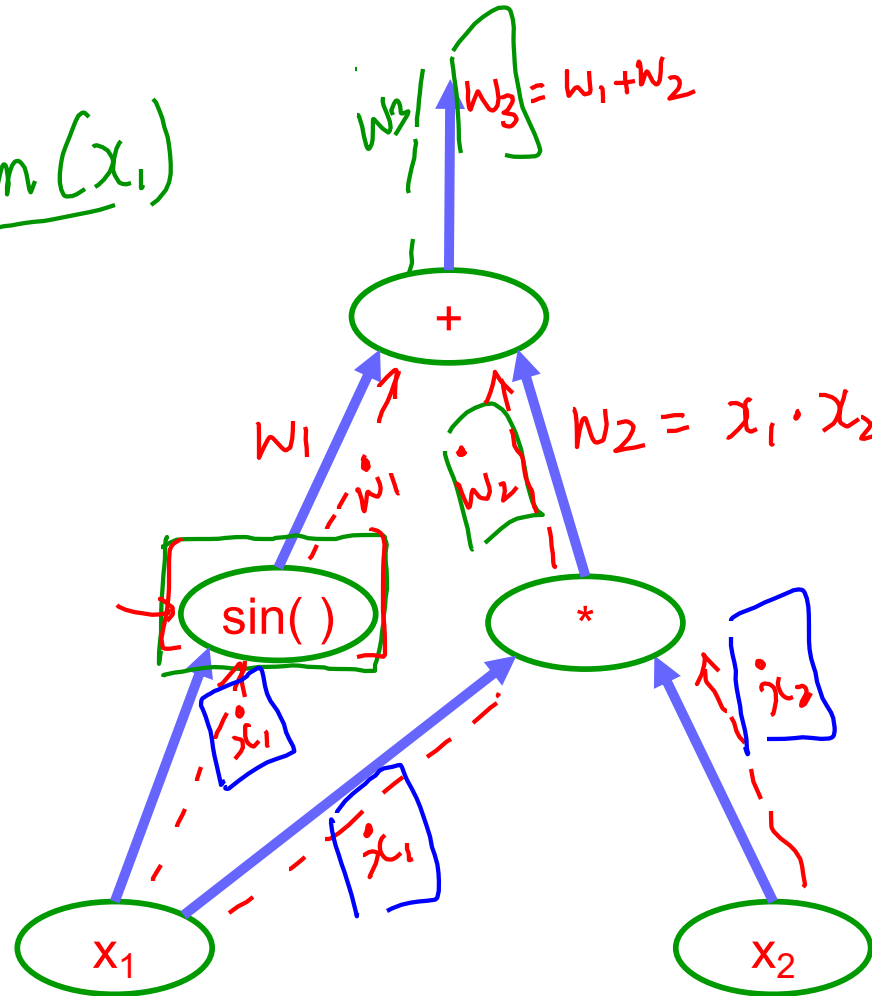
$\dfrac{\partial f}{\partial x_2}$

$= x_1$



$x_1$

$\partial \cdot$

$x_1$

$\dfrac{\partial x_1}{\partial x_2}$  $0$

sin( )  $*$

$\dfrac{\partial x_2}{\partial x_2} = 1$

$x_1$  $x_2$

# Example: Forward mode AD

Goal:

$$f(x_1, x_2) = \overbrace{x_1 x_2}^{w_2} + \overbrace{\sin(x_1)}^{w_1}$$

$$\frac{\partial f}{\partial x} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

$$\frac{\partial x_1}{\partial a} = \dot{x}_1$$

$$\frac{\partial x_2}{\partial a} = \quad \exists$$

$$W_1 = \sin(x_1)$$

Input

$$\boxed{\frac{\partial W_1}{\partial a}} = \boxed{\frac{\partial W_1}{\partial x_1}} \cdot \boxed{\frac{\partial x_1}{\partial a}}$$

$$\cos(x_1) \cdot \dot{x}_1$$

$$\frac{\partial W_2}{\partial a}$$

$$\boxed{\frac{\partial W_3}{\partial a}}$$

$$\frac{\partial f}{\partial a}$$

$W_3 = W_1 + W_2$

$W_1$

$\dot{w}_1$

$W_2 = x_1 \cdot x_2$

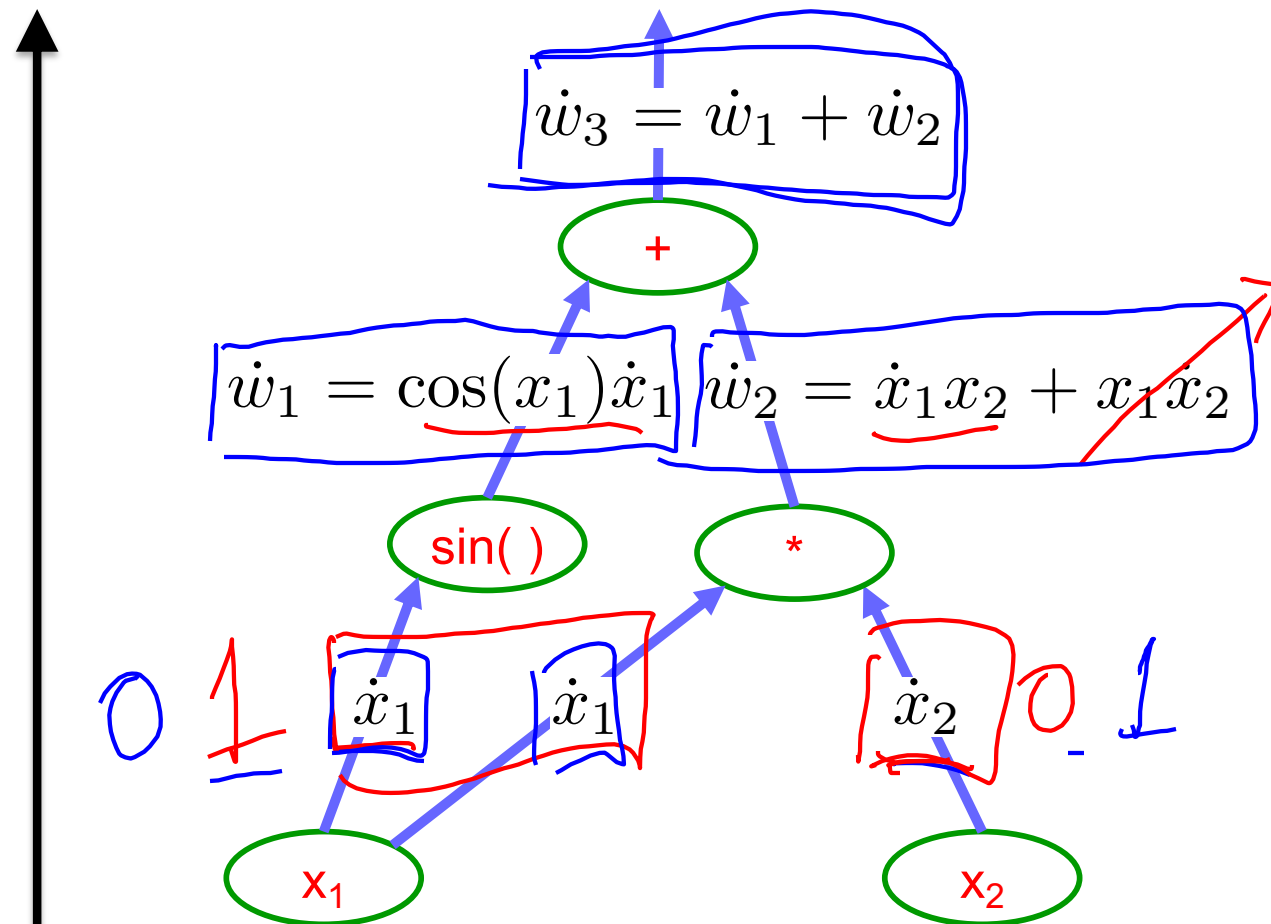$\dot{w}_2$

sin( )

$\dot{x}_1$

$*$

$\dot{x}_2$

$\dot{x}_1$

$x_1$

$x_2$

$$\frac{\partial W_2}{\partial a} = x_1 \frac{\partial x_2}{\partial a} + x_2 \frac{\partial x_1}{\partial a}$$

$$= x_1 \dot{x}_2 + x_2 \dot{x}_1$$

$$a \in \{x_1, x_2\}$$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

$$\frac{\partial f}{\partial x} = \left[ \frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right]$$

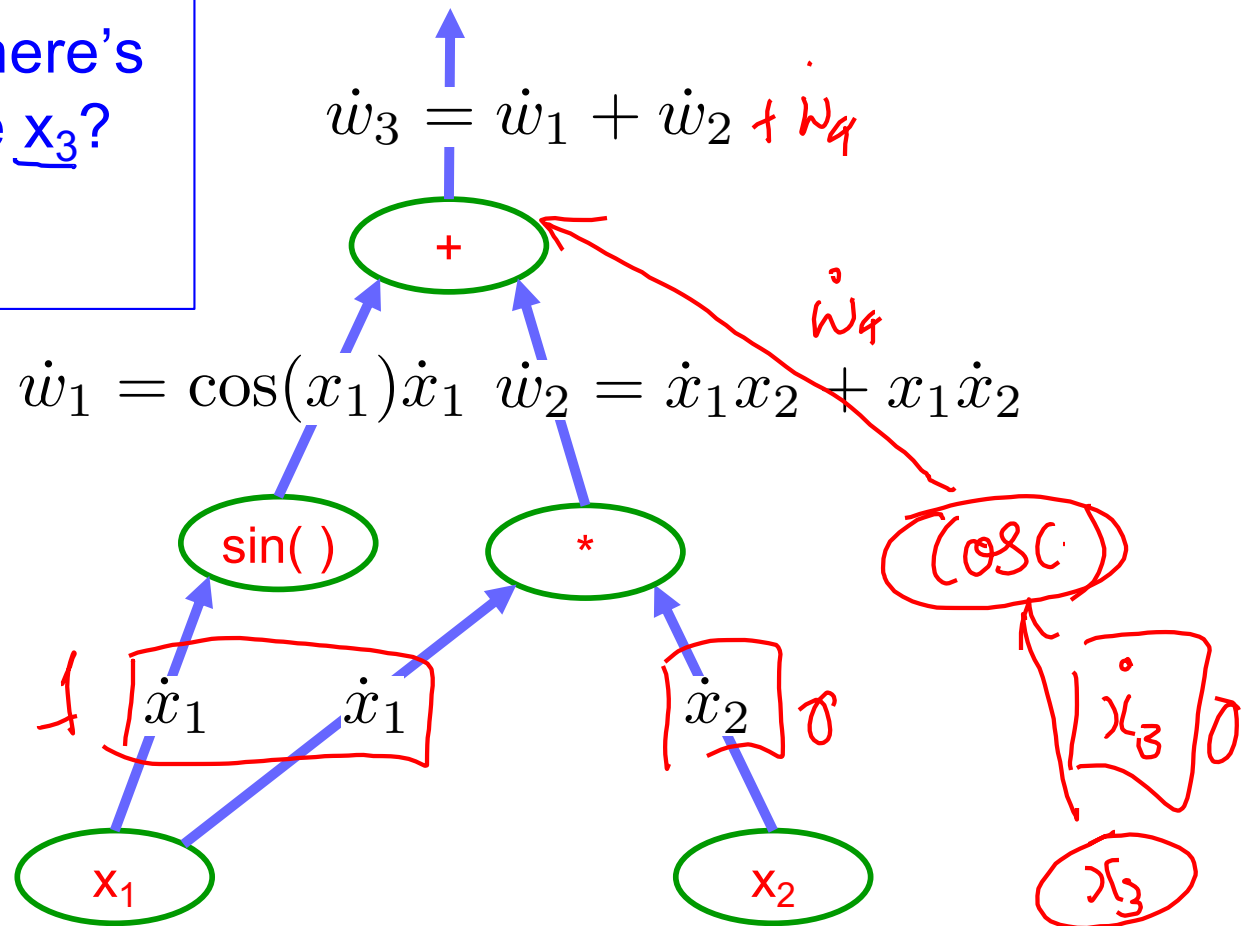$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$+$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 x_2$$

sin( )     *

$$\dot{x}_1 \quad \dot{x}_1 \quad \dot{x}_2$$

0  1                    0  1

$x_1$       $x_2$

$$\dot{x}_1 = \frac{\partial x_1}{\partial a}$$

$$\dot{w}_1 = \frac{\partial w_1}{\partial a}$$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1) + \cos(x_3)$$

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2 + \dot{w}_4$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \qquad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$
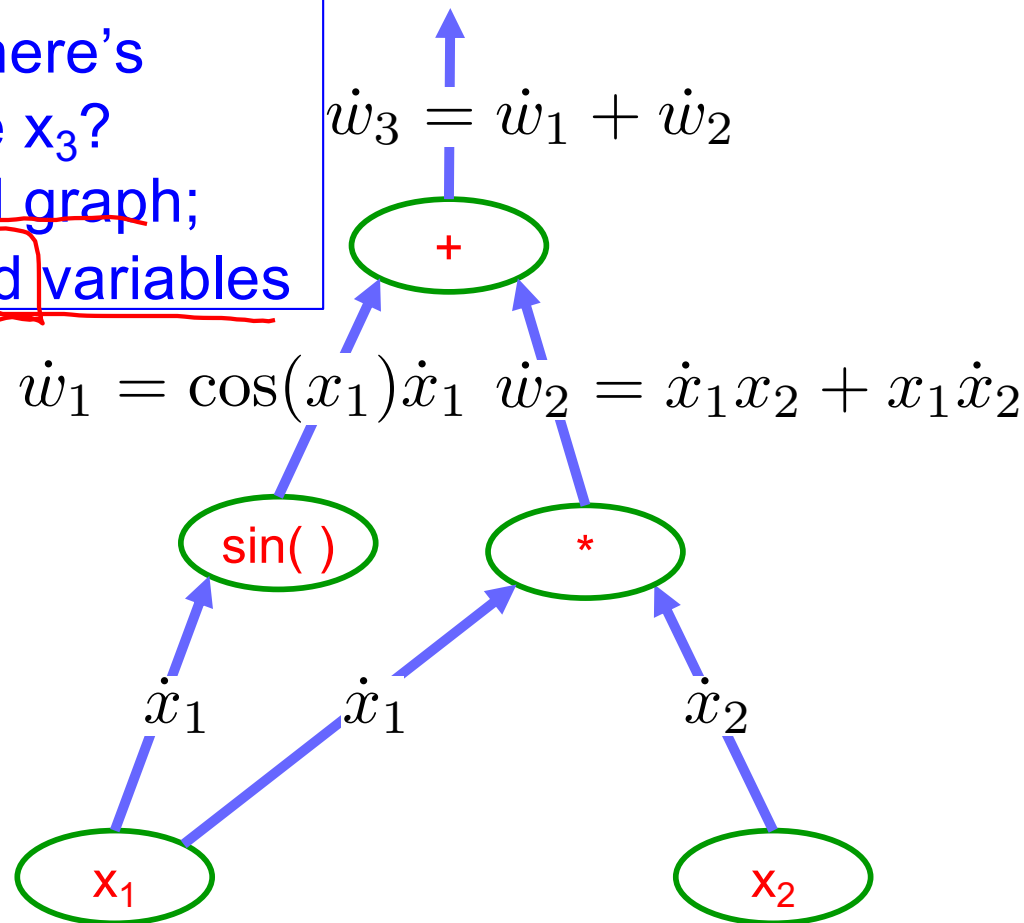
$$\dot{w}_4$$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$
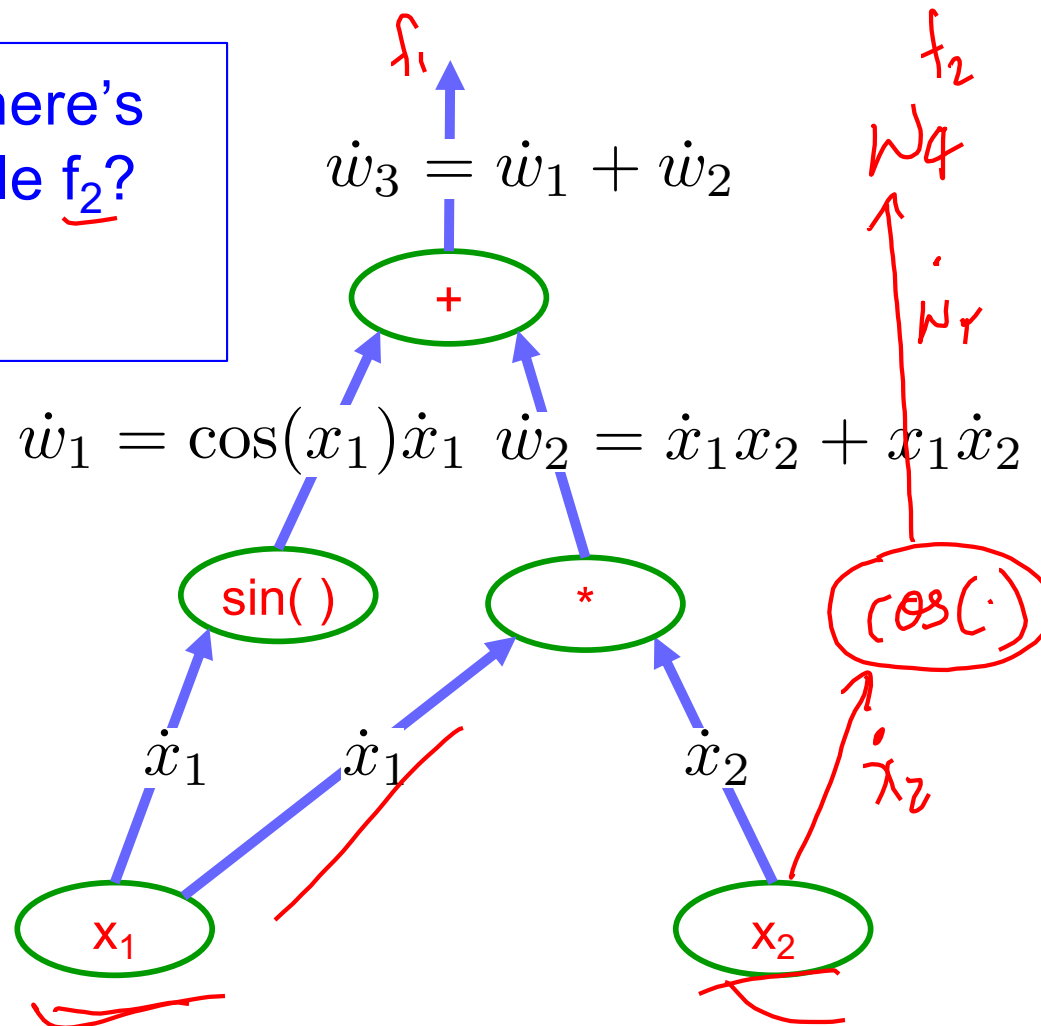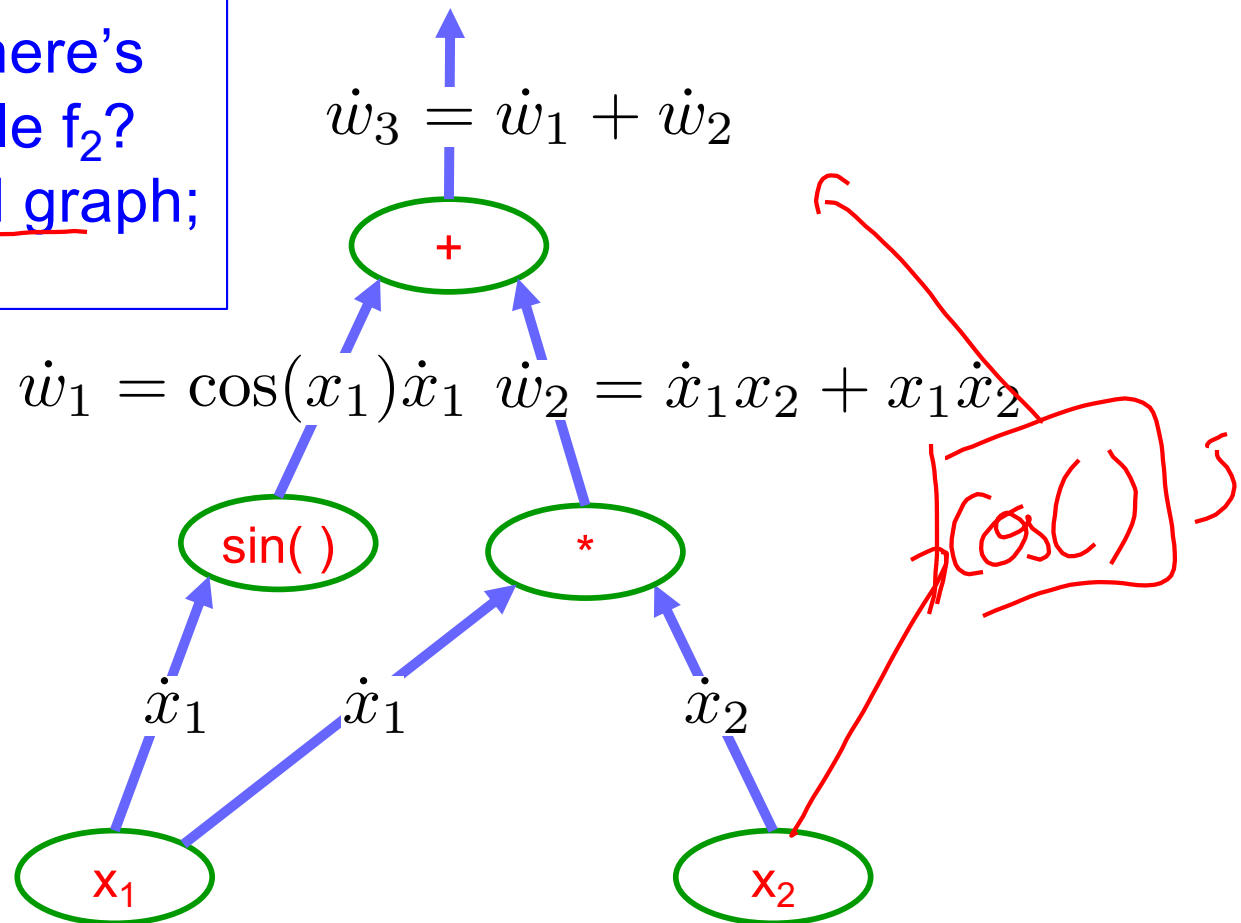
Q: What happens if there's another input variable $x_3$?
A: more sophisticated graph; d "forward props" for d variables

$$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$$

$$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$$

+

sin( )          *

$\dot{x}_1$        $\dot{x}_1$        $\dot{x}_2$

$x_1$                          $x_2$

# Example: Forward mode AD

$$f_1(x_1, x_2) = x_1 x_2 + \sin(x_1) \quad f_2 = \cos(x_2)$$

$f_1$

$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

$f_2$

$\dot{w}_4$

$+$

$\dot{w}_r$

$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

sin( )     *     $\cos(\cdot)$

$\dot{x}_1$    $\dot{x}_1$      $\dot{x}_2$    $\dot{x}_2$

$x_1$               $x_2$

# Example: Forward mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

Q: What happens if there's another output variable $f_2$?
A: more sophisticated graph; single "forward prop"

2

for $a \in \{x_1, x_2\}$

$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

$+$

$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

$\cos()$

sin( )   *

$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$

$x_1$   $x_2$

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



$w_3 = w_1 + w_2$

$\overline{w_1} = \dfrac{\partial f}{\partial w_1} = \dfrac{\partial f}{\partial w_3} \cdot \dfrac{\partial w_3}{\partial w_1}$

$= \overline{w_3} = 1$

$\overline{w_3} = \dfrac{\partial f}{\partial w_3} = \dfrac{\partial w_3}{\partial w_3} = 1$

$\overline{w_2} = \dfrac{\partial f}{\partial w_2} \qquad \overline{w_1} = \dfrac{\partial f}{w_1}$

$w_1 = \sin(x_1)$

$\overline{x_1} = \dfrac{\partial f}{\partial x_1} = \dfrac{\partial f}{\partial w_1} \cdot \dfrac{\partial w_1}{\partial x_1}$

$\cos(x_1)$

$\overline{x_1} = \dfrac{\partial f}{\partial x_1} \qquad \overline{x_2} = \dfrac{\partial f}{\partial x_2}$
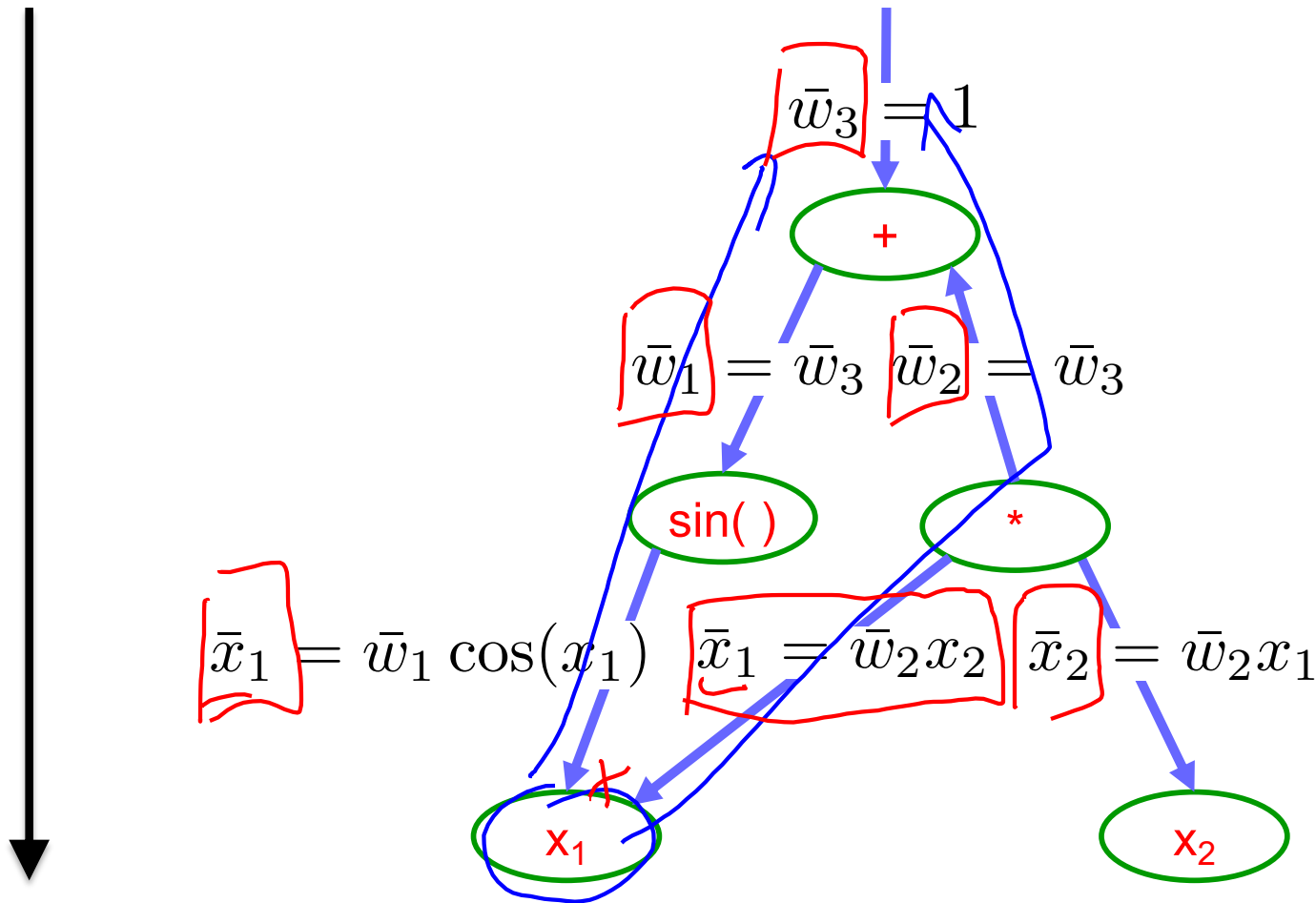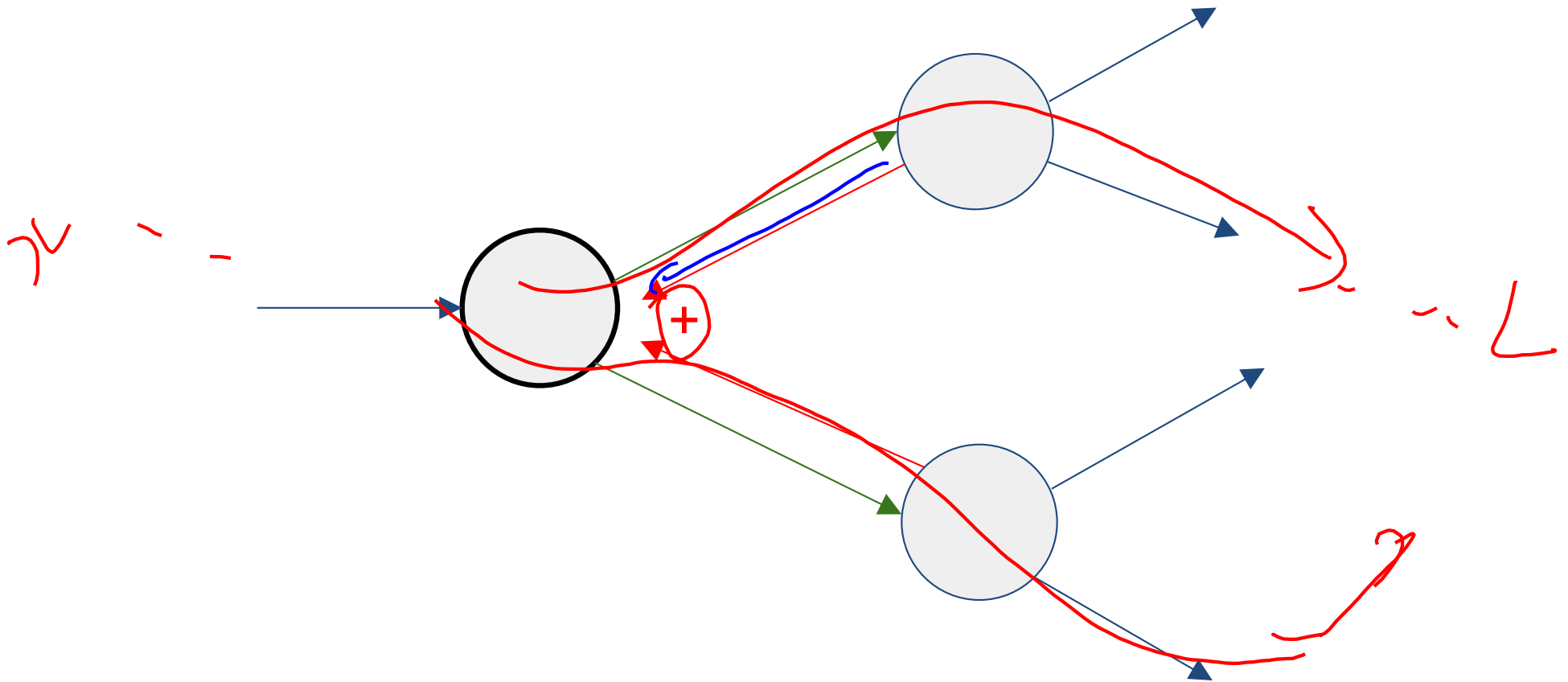
(C) Dhruv Batra

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

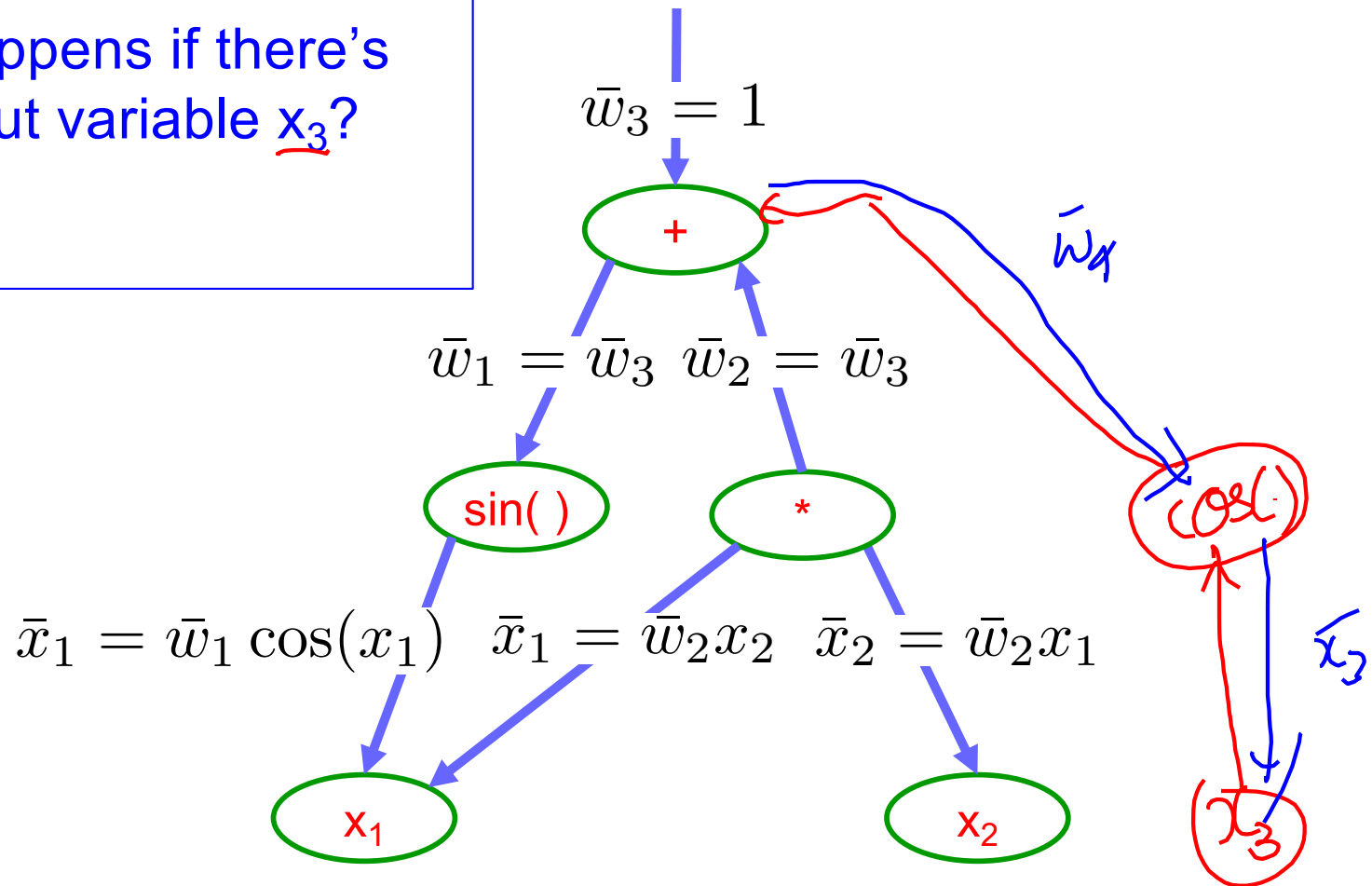# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



$\bar{w}_3 = 1$

$+$

$\bar{w}_1 = \bar{w}_3$ $\quad$ $\bar{w}_2 = \bar{w}_3$

$\sin(\ )$ $\qquad$ $*$

$\bar{x}_1 = \bar{w}_1 \cos(x_1)$ $\quad$ $\bar{x}_1 = \bar{w}_2 x_2$ $\quad$ $\bar{x}_2 = \bar{w}_2 x_1$

$x_1$ $\qquad$ $x_2$

# Gradients add at branches

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1) + \cos(x_3)$$
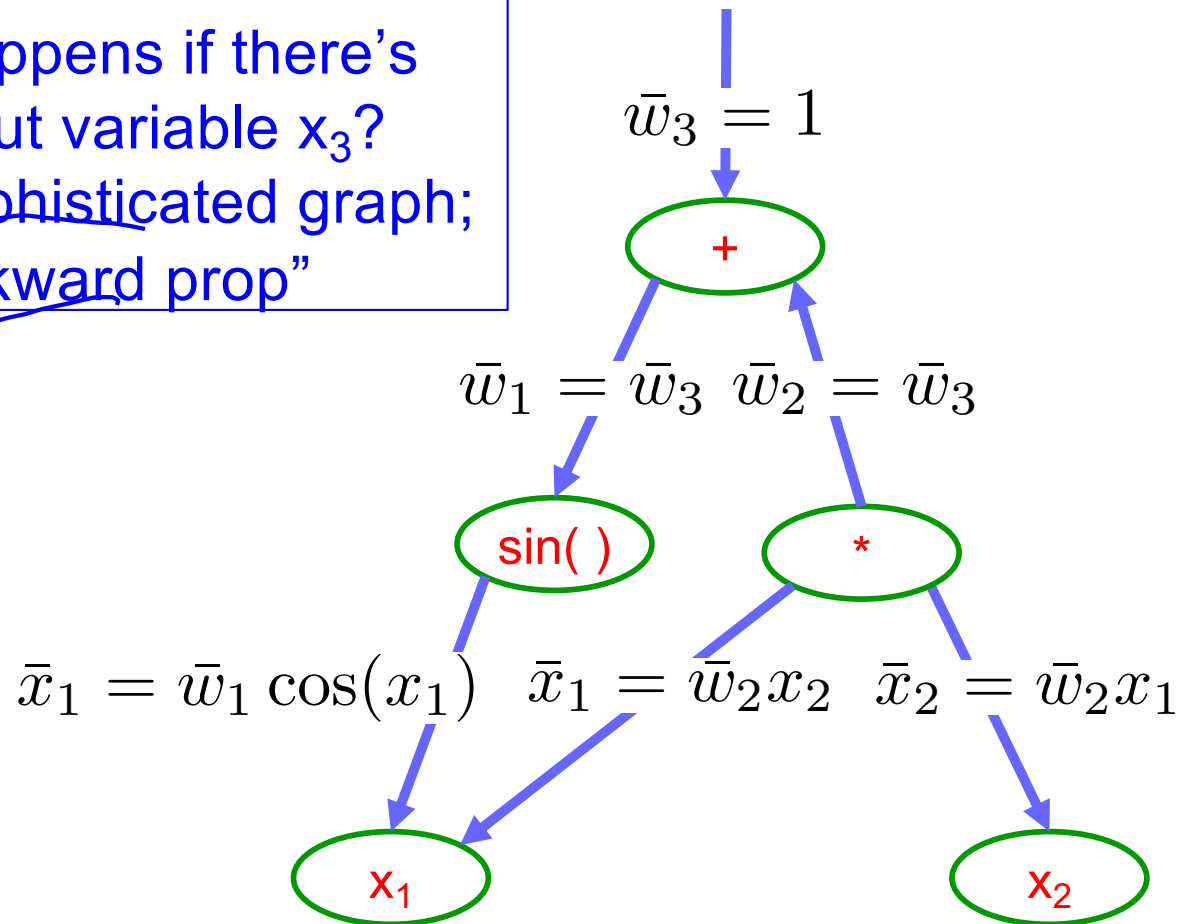
Q: What happens if there's another input variable $x_3$?

$\bar{w}_3 = 1$

$+$

$\bar{w}_1 = \bar{w}_3$  $\bar{w}_2 = \bar{w}_3$

$\bar{w}_4$

$\sin( )$      $*$      $\cos( )$

$\bar{x}_1 = \bar{w}_1 \cos(x_1)$  $\bar{x}_1 = \bar{w}_2 x_2$  $\bar{x}_2 = \bar{w}_2 x_1$

$x_3$

$x_1$      $x_2$      $x_3$
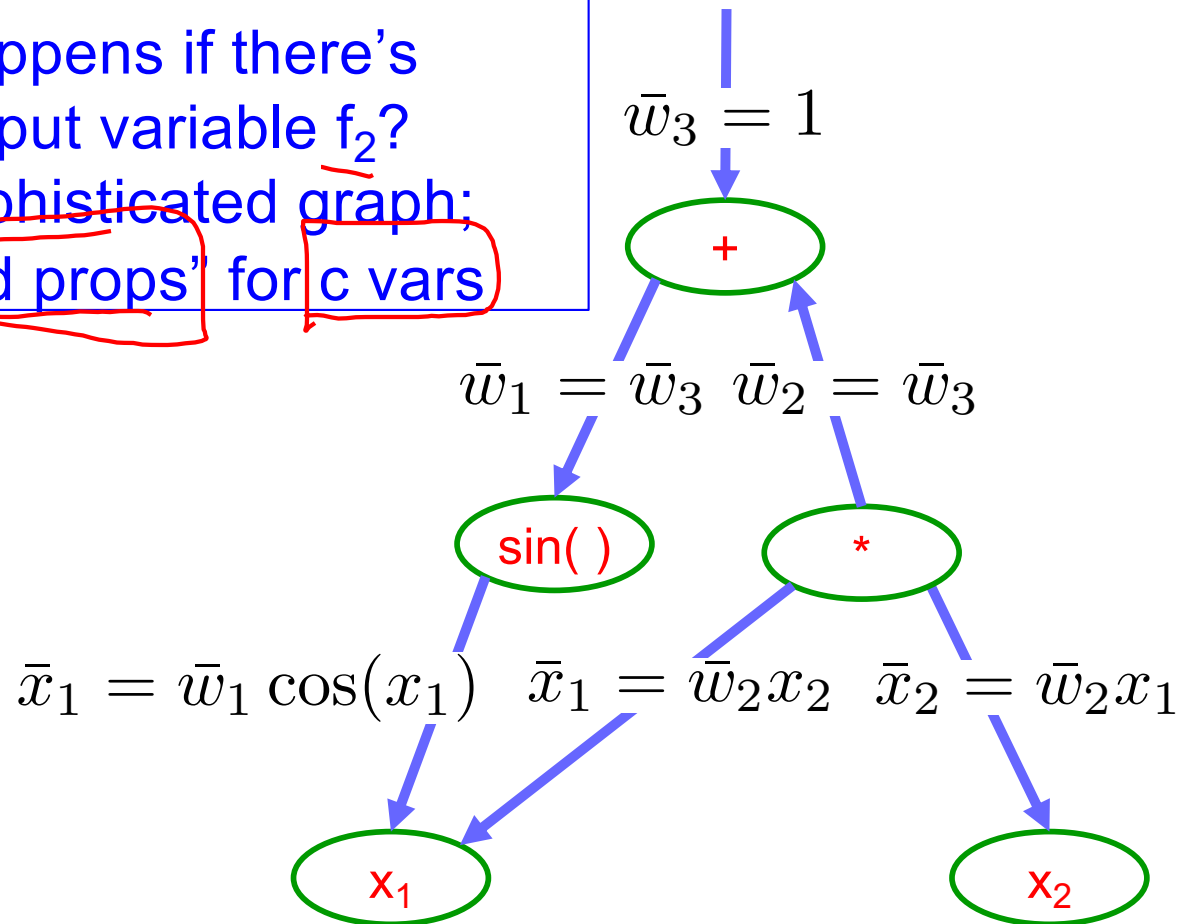
# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

Q: What happens if there's another input variable x₃?
A: more sophisticated graph; single "backward prop"

$\bar{w}_3 = 1$

$+$

$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$

$\sin(\ )$ $*$

$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$

$x_1$ $x_2$

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

$f_2 = \cos(x_2)$

Q: What happens if there's another output variable $f_2$?

For an output vars

$\bar{w}_3 = 1$ $\{1,0\}$

$W_4 = f_2$

$\bar{w}_4 = \{0,1\}$

+

$\bar{w}_1 = \bar{w}_3$ $\bar{w}_2 = \bar{w}_3$

sin( )

*

cos( )

$\bar{x}_1 = \bar{w}_1 \cos(x_1)$ $\bar{x}_1 = \bar{w}_2 x_2$ $\bar{x}_2 = \bar{w}_2 x_1$

$x_1$

$x_2$

# Example: Reverse mode AD

$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$

*l . . . . fc*

Q: What happens if there's
another output variable f_2?
A: more sophisticated graph;
c "backward props" for c vars



$\bar{w}_3 = 1$

$+$

$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$

sin( )   *

$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$

x_1   x_2

# Forward mode vs Reverse Mode

- x → Graph → L
- Intuition of Jacobian

# Forward mode vs Reverse Mode

- What are the differences?

- Which one is faster to compute?
  - Forward or backward? Is $c > d$ or $c < d$

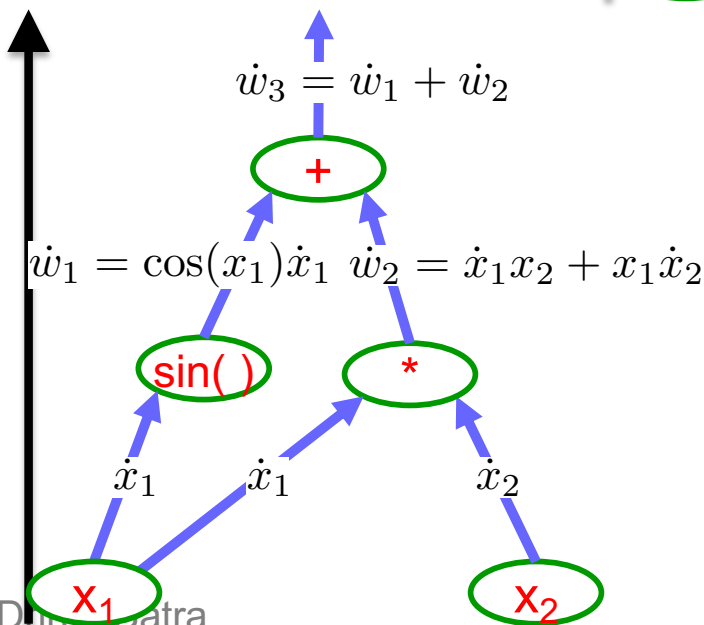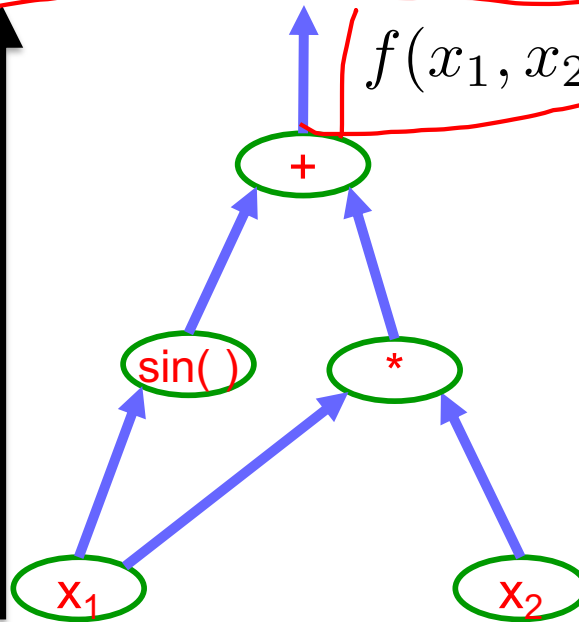# Forward mode vs Reverse Mode

- What are the differences?

- Which one is faster to compute?
  - Forward or backward?

- Which one is more memory efficient (less storage)?
  - Forward or backward?
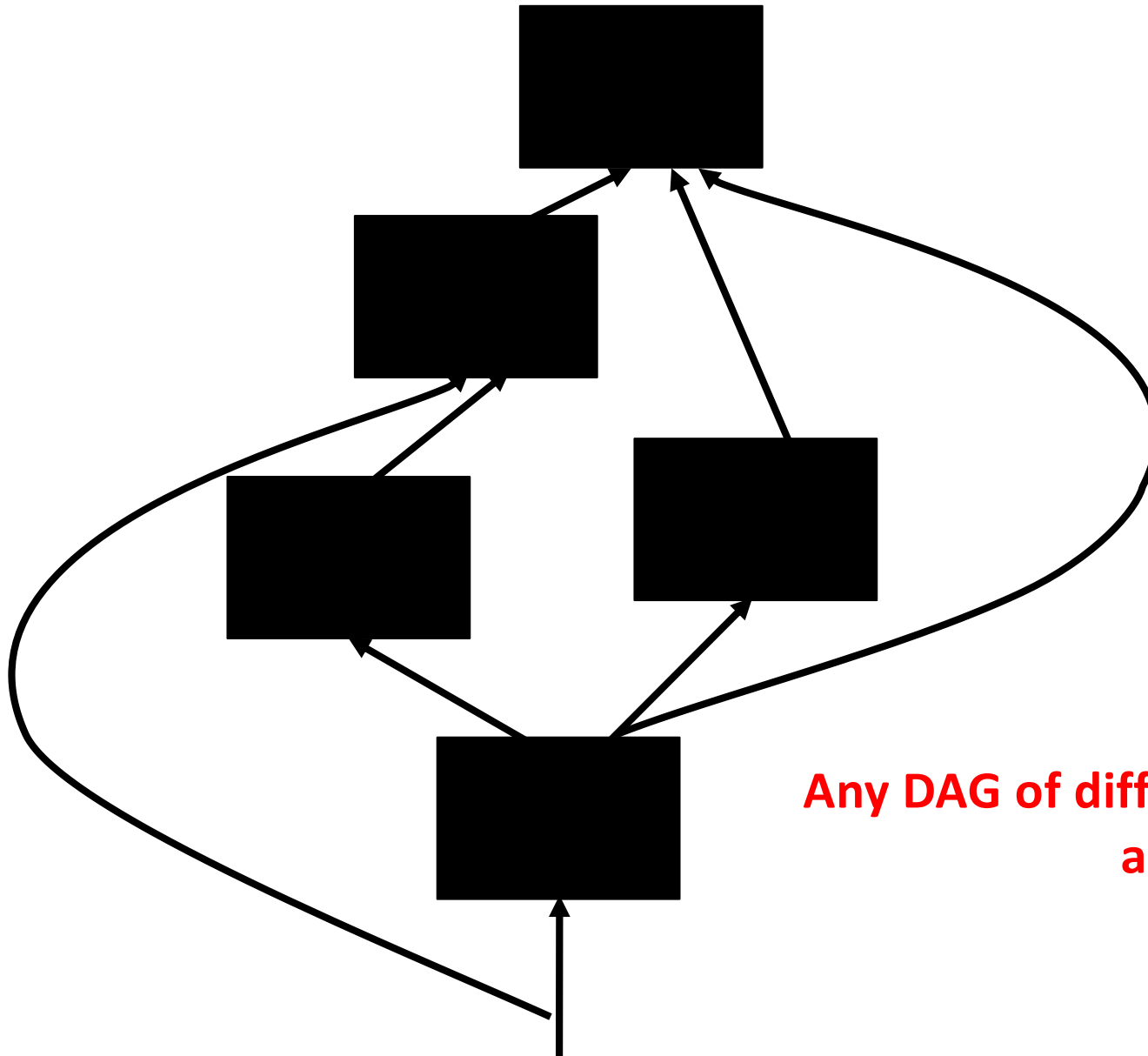
# Forward Pass vs
# Forward mode AD vs Reverse Mode AD
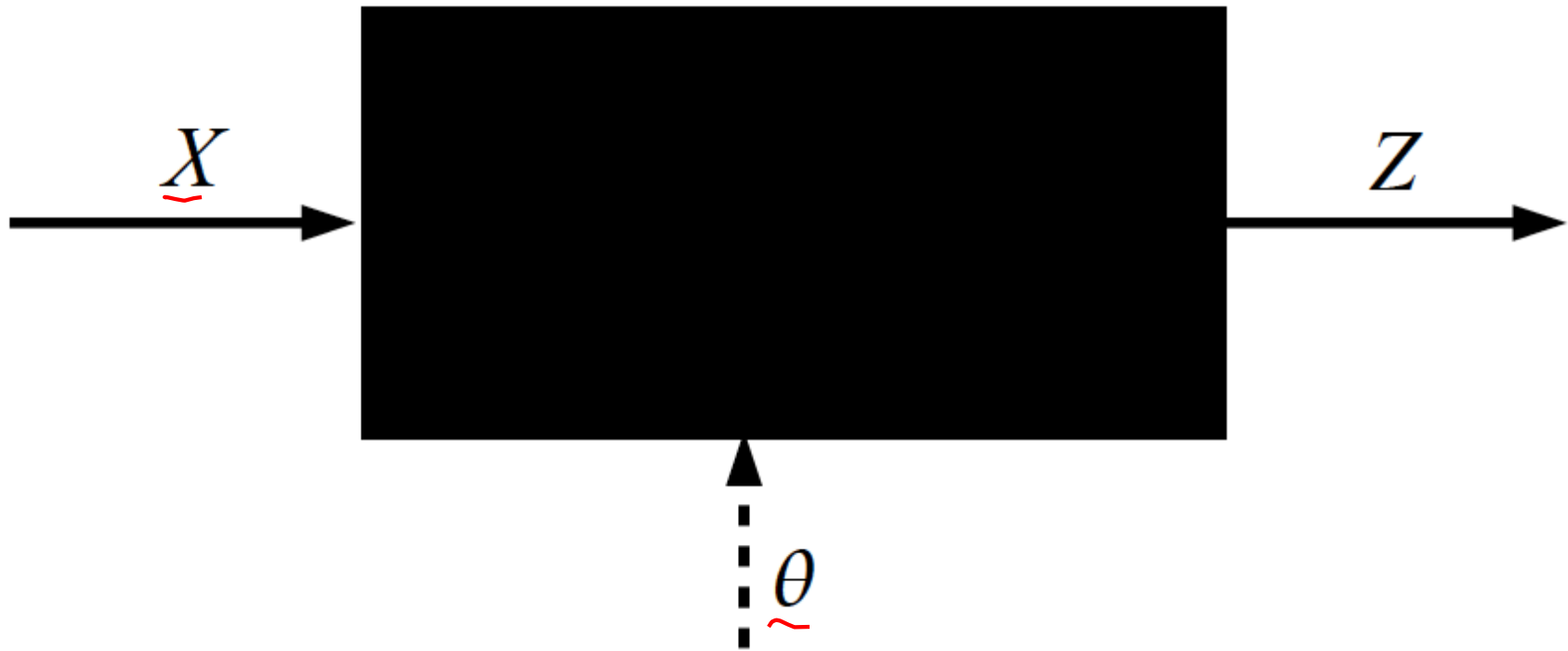
$$f(x_1, x_2) = x_1 x_2 + \sin(x_1)$$



$\dot{w}_3 = \dot{w}_1 + \dot{w}_2$

$\dot{w}_1 = \cos(x_1)\dot{x}_1 \quad \dot{w}_2 = \dot{x}_1 x_2 + x_1 \dot{x}_2$

$\dot{x}_1 \qquad \dot{x}_1 \qquad \dot{x}_2$

$\bar{w}_3 = 1$

$\bar{w}_1 = \bar{w}_3 \quad \bar{w}_2 = \bar{w}_3$

$\bar{x}_1 = \bar{w}_1 \cos(x_1) \quad \bar{x}_1 = \bar{w}_2 x_2 \quad \bar{x}_2 = \bar{w}_2 x_1$

43

# Plan for Today

- Automatic Differentiation
  - (Finish) Forward mode vs Reverse mode AD ✓
  - Backprop
  - Patterns in backprop
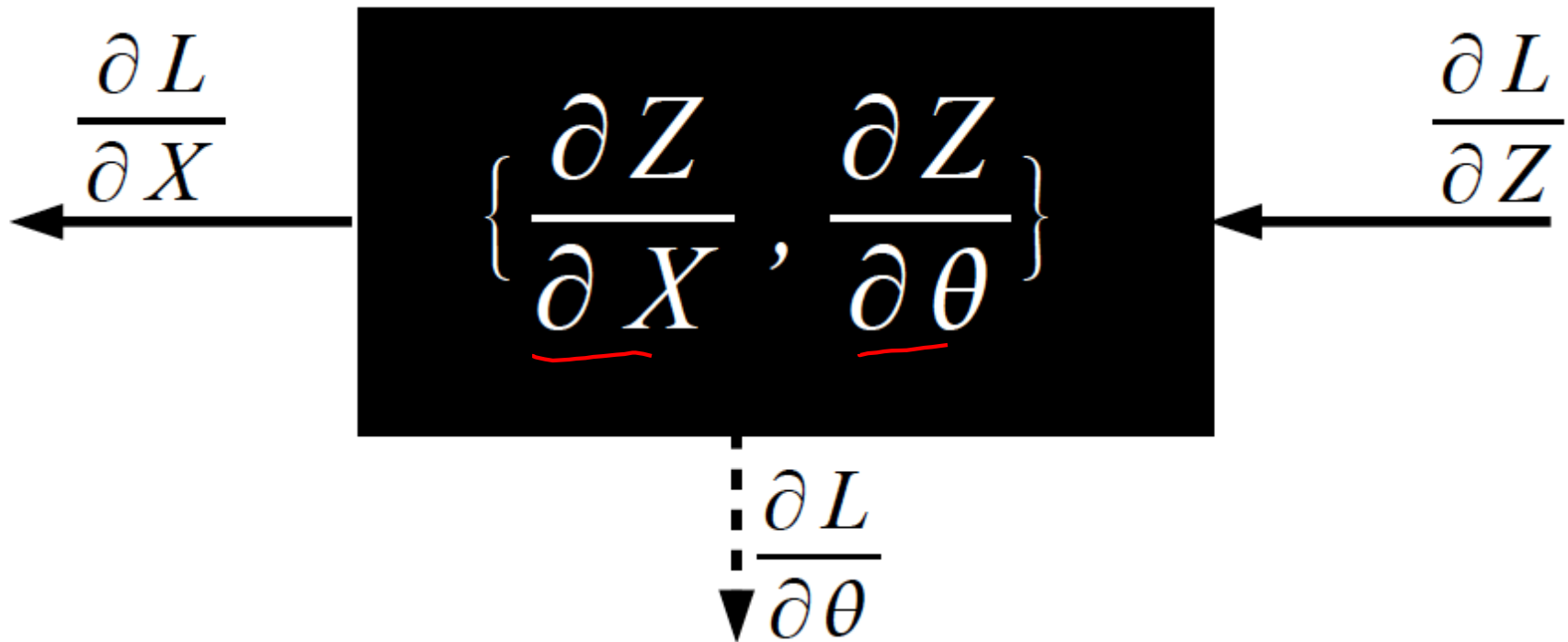  - Jacobians in FC+ReLU NNs

# Computational Graph



**Any DAG of differentiable modules is allowed!**
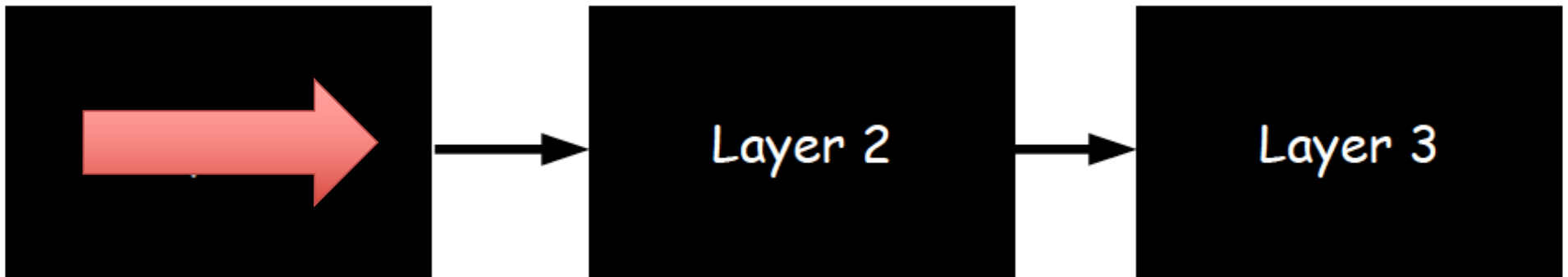
# Key Computation: Forward-Prop

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Key Computation: Back-Prop



$$\frac{\partial L}{\partial X} \qquad \left\{ \frac{\partial Z}{\partial X}, \frac{\partial Z}{\partial \theta} \right\} \qquad \frac{\partial L}{\partial Z}$$
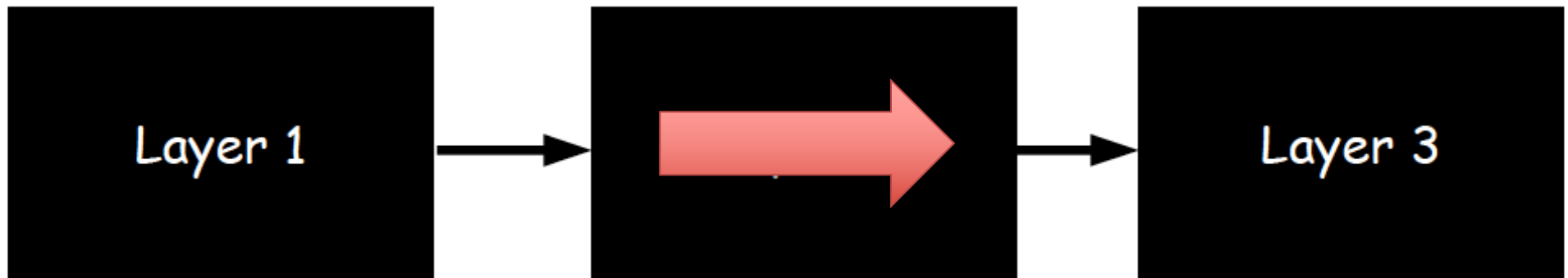
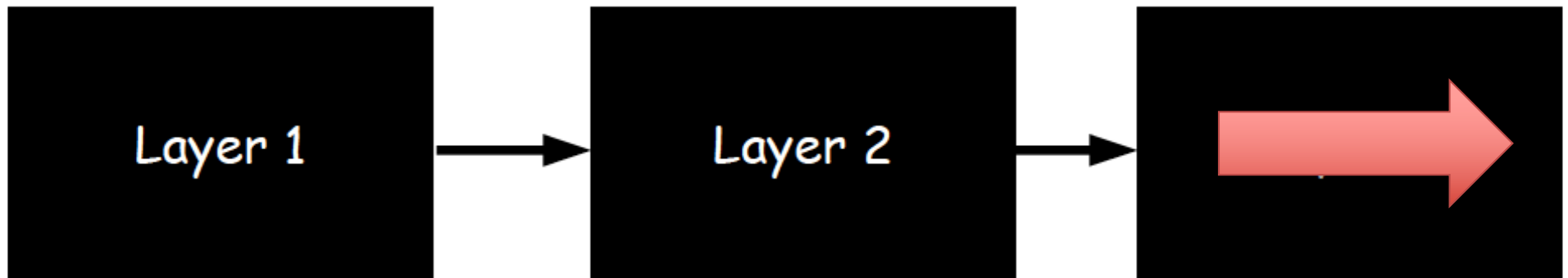$$\frac{\partial L}{\partial \theta}$$

# Neural Network Training

- Step 1: Compute Loss on mini-batch    [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch        [F-Pass]
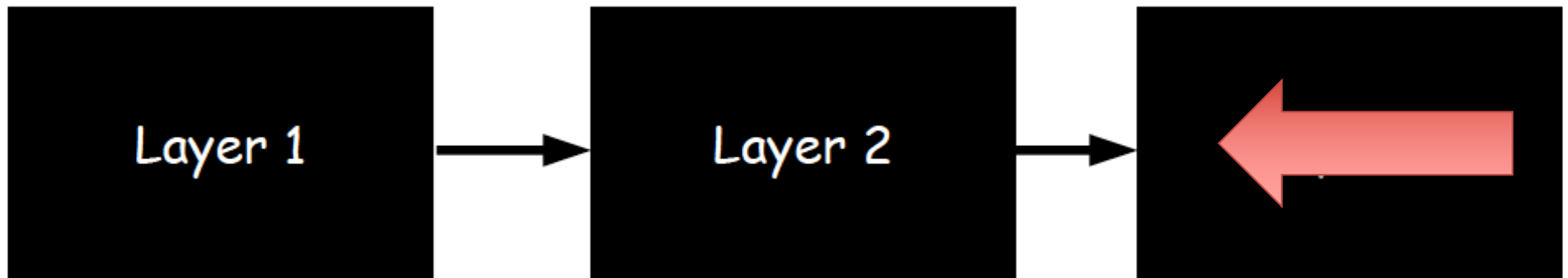
# Neural Network Training

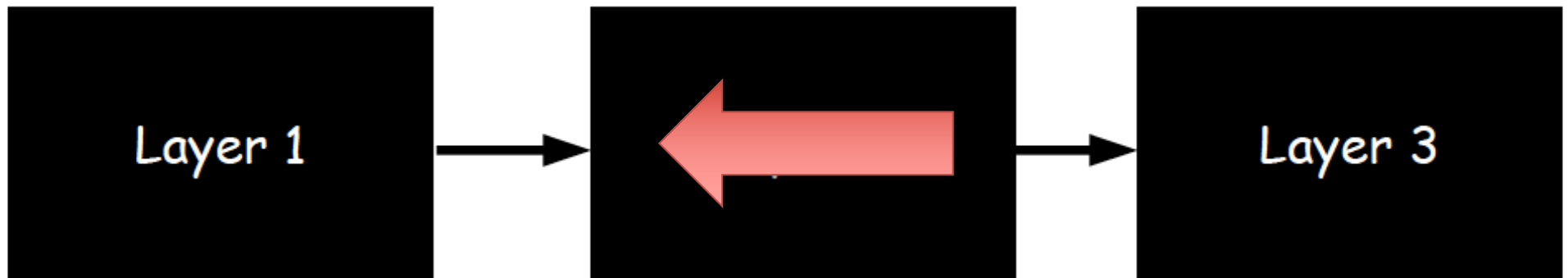- Step 1: Compute Loss on mini-batch          [F-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch        [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch      [F-Pass]
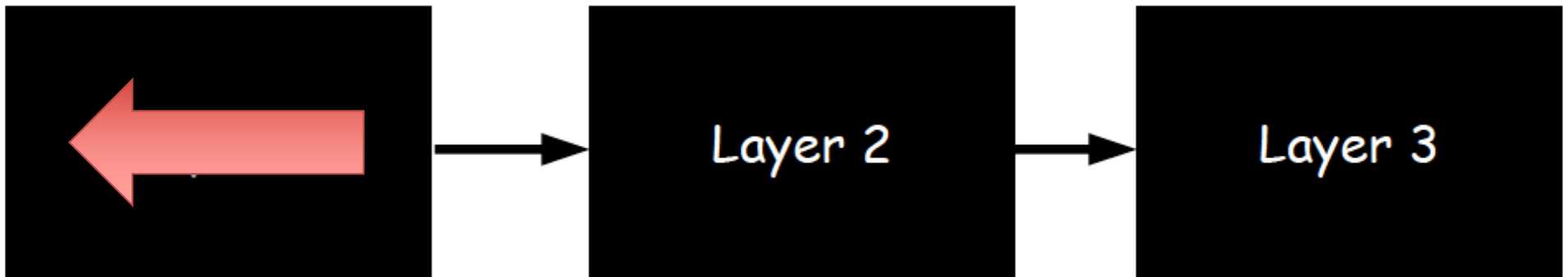- Step 2: Compute gradients wrt parameters   [B-Pass]
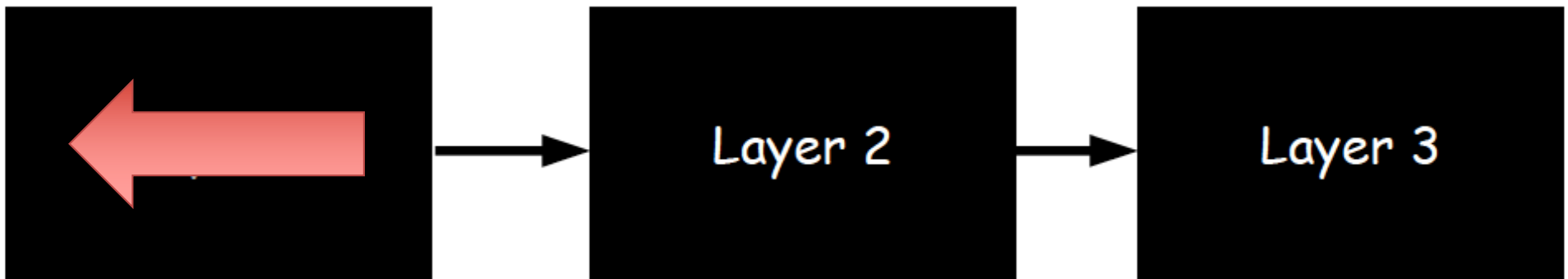
# Neural Network Training

- Step 1: Compute Loss on mini-batch    [F-Pass]
- Step 2: Compute gradients wrt parameters   [B-Pass]

# Neural Network Training

- Step 1: Compute Loss on mini-batch       [F-Pass]
- Step 2: Compute gradients wrt parameters    [B-Pass]
- Step 3: Use gradient to update parameters



$$\theta \leftarrow \theta - \eta \frac{dL}{d\theta}$$

$$\mathbf{x} \longrightarrow \boxed{f_1} \xrightarrow{\ \mathbf{h}^{(1)}\ } \boxed{f_2} \xrightarrow{\ \mathbf{h}^{(2)}\ } \mathbf{h}^{(L-1)} \cdots \longrightarrow \boxed{f_L} \longrightarrow L \in \mathbb{R}$$

$\mathbf{w}_1 \qquad \mathbf{w}_2 \qquad \dfrac{\partial L}{\partial w_2} \qquad \mathbf{w}_L$

$\dfrac{\partial L}{\partial w_1}$