

# Autodiff

*The algorithm that will upend the world (maybe)*

Lim Soon Wei, Daniel

# Norming

*(verb) forming shared “norms”: expectations, styles, comfort*

Questions are  
welcome!

Ask directly in the chat,  
raise a “hand”, or just  
unmute.

Enjoy the story

Don't worry about taking  
notes; you can download  
these slides and demo  
code later.

No pressure

I won't call on you unless  
you raise a “hand”

*Switch on your camera only if you are comfortable!*

**Have you used:**

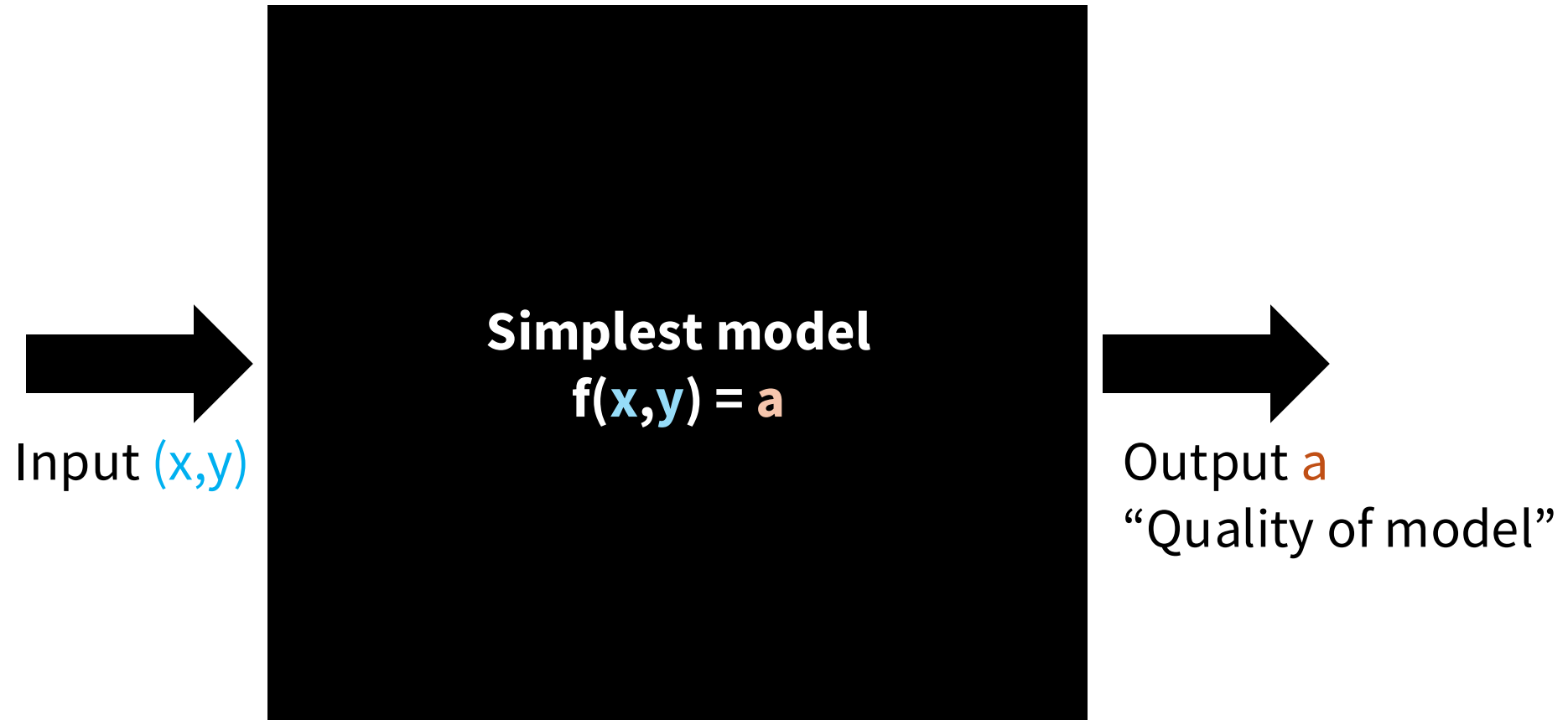
 Gemini

 deepseek

 ChatGPT

**Machine learning model**

10101011001100101010100101010101  
01001010101010100101010100110101  
011001100101010100101010101001  
01010101010010101010011010101100  
11001010101001010101010100101010  
10101001010101001101010110011001  
01010100101010101010010101010101  
**Machine learning model**  
00101010100110101011001100101010  
100101010101001010101010100101  
01010011010101100110010101010010  
10101010100101010101010010101010  
01101010110011001010101001010101  
01010010101010101001010101001101  
01011001100101010100101010101010

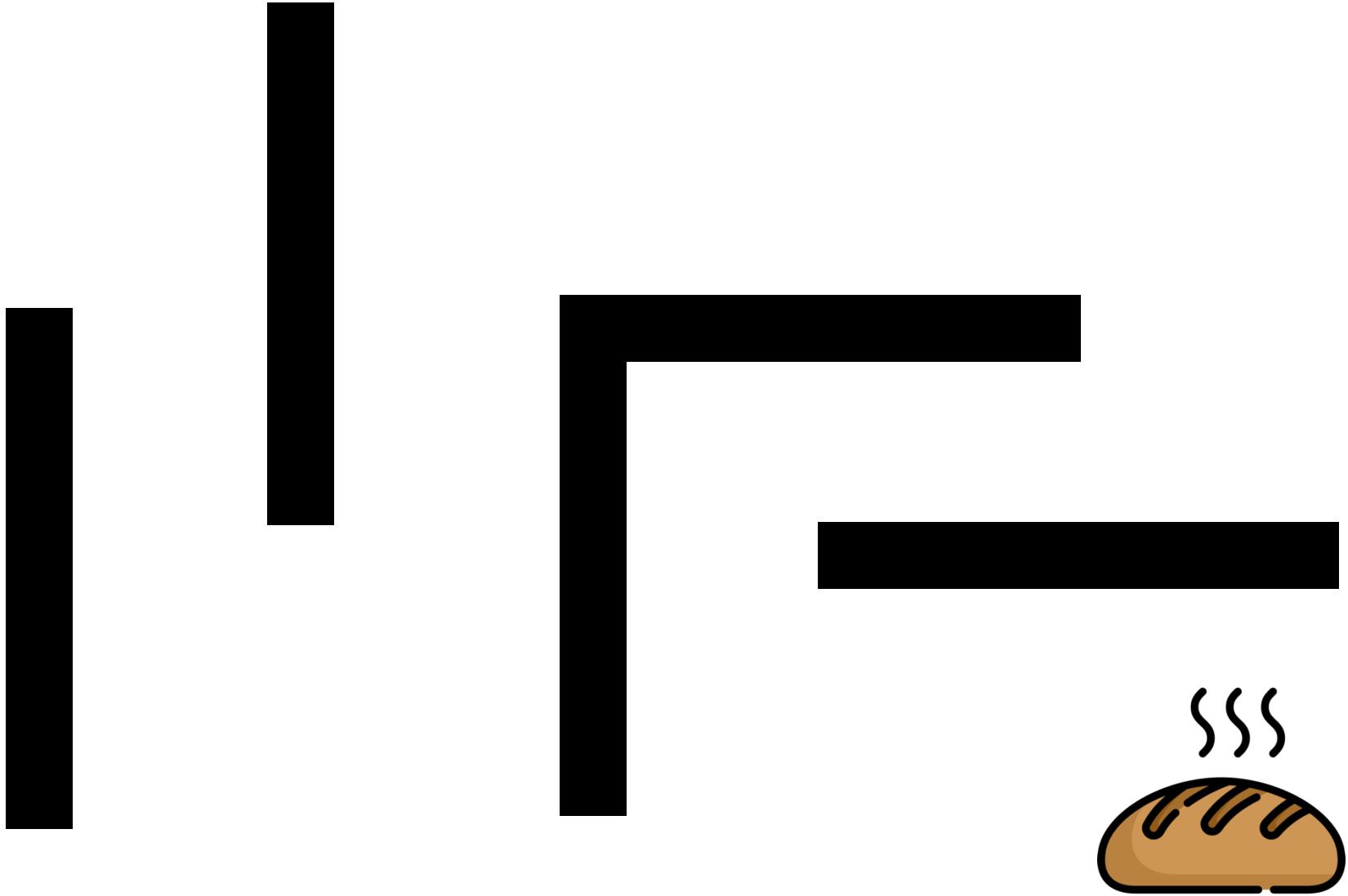


Optimization problem: How to pick  $(x,y)$  to minimize/maximize  $a$ ?



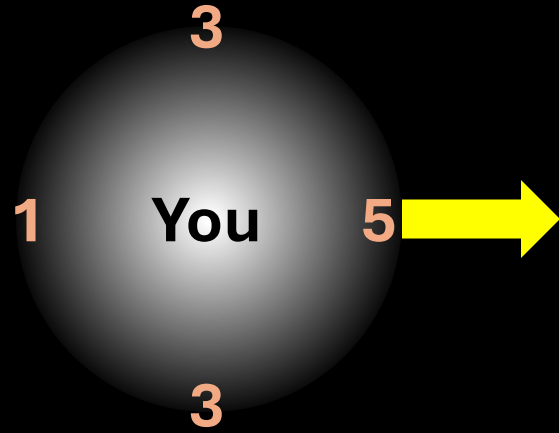


You

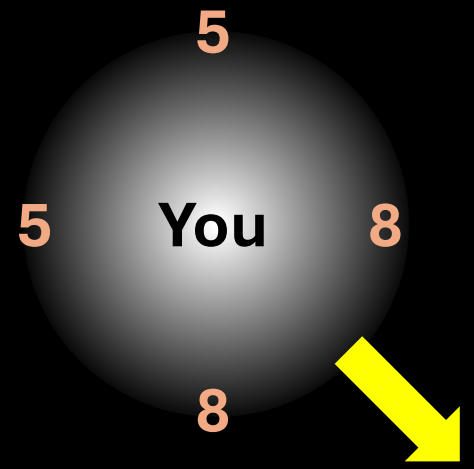




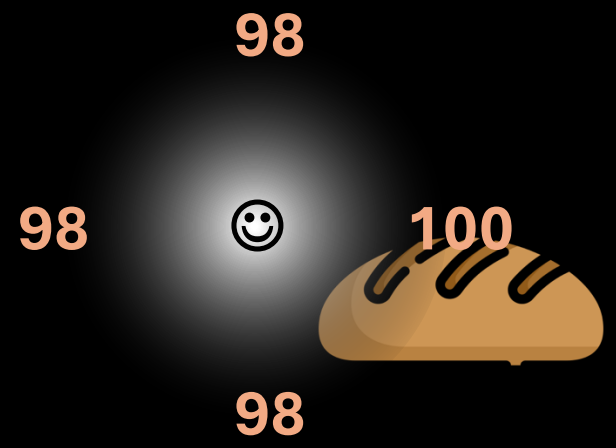
“Aroma” score



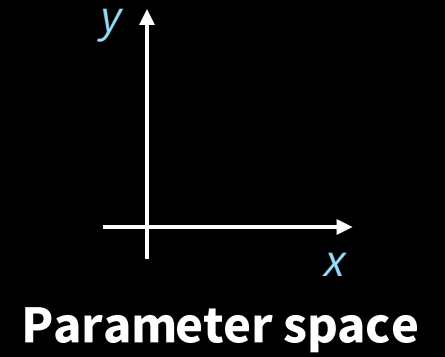
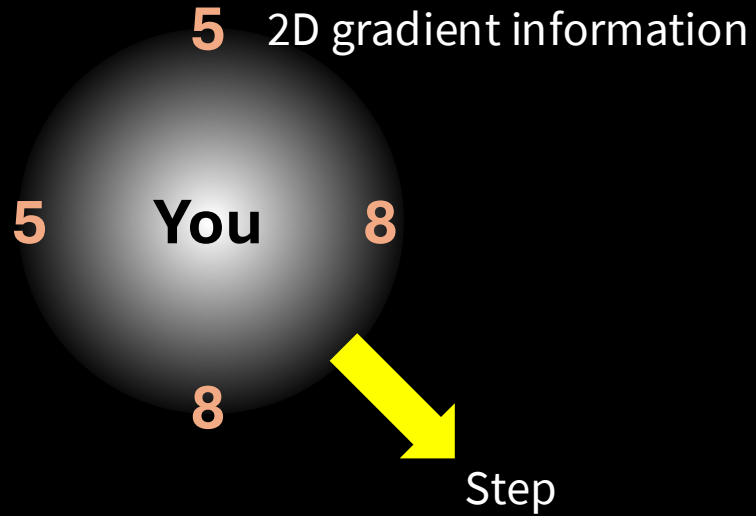
“Aroma” score



“Aroma” score



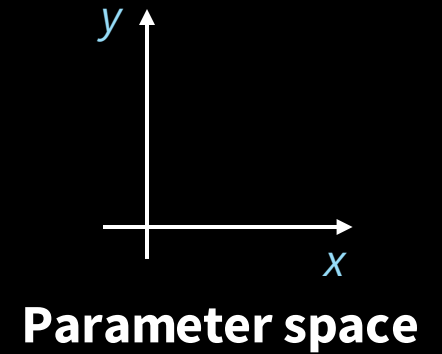
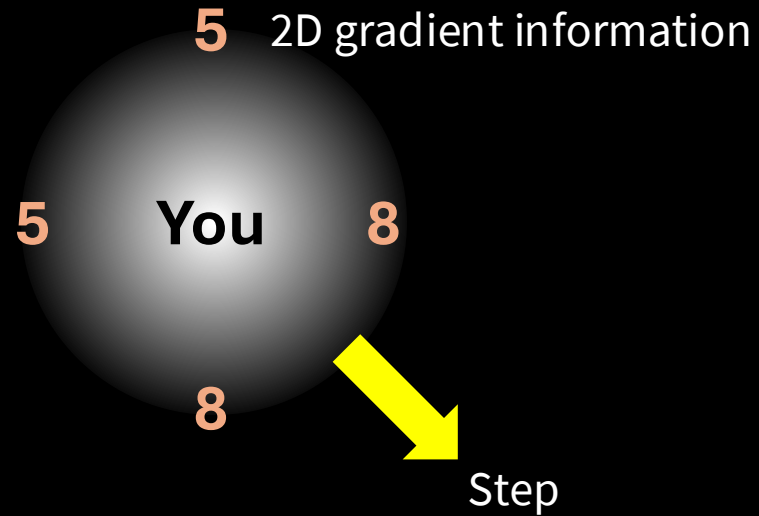
# Gradient descent



## Algorithm

1. Compute gradient
2. Move a step in direction of greatest increase/decrease
3. Good enough? If not, go back to Step 1.

# Gradient descent (mathematically)



## Algorithm

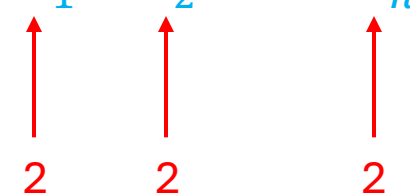
1. Compute gradient  $\nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$  at current position  $(x_0, y_0)$
2. Move a step:  $(x_0, y_0) \mapsto (x_0, y_0) - L \nabla f$ ,  $L =$  step size
3. Check for convergence (e.g., is  $\nabla f$  small?) . If not, go back to Step 1.

# How many steps does it take in finite differences?

$$\text{In 2D: } \nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

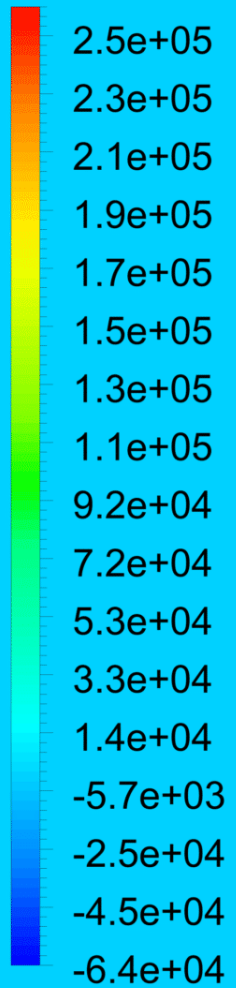
$$\text{In 3D: } \nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right)$$

$$\text{In nD: } \nabla f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

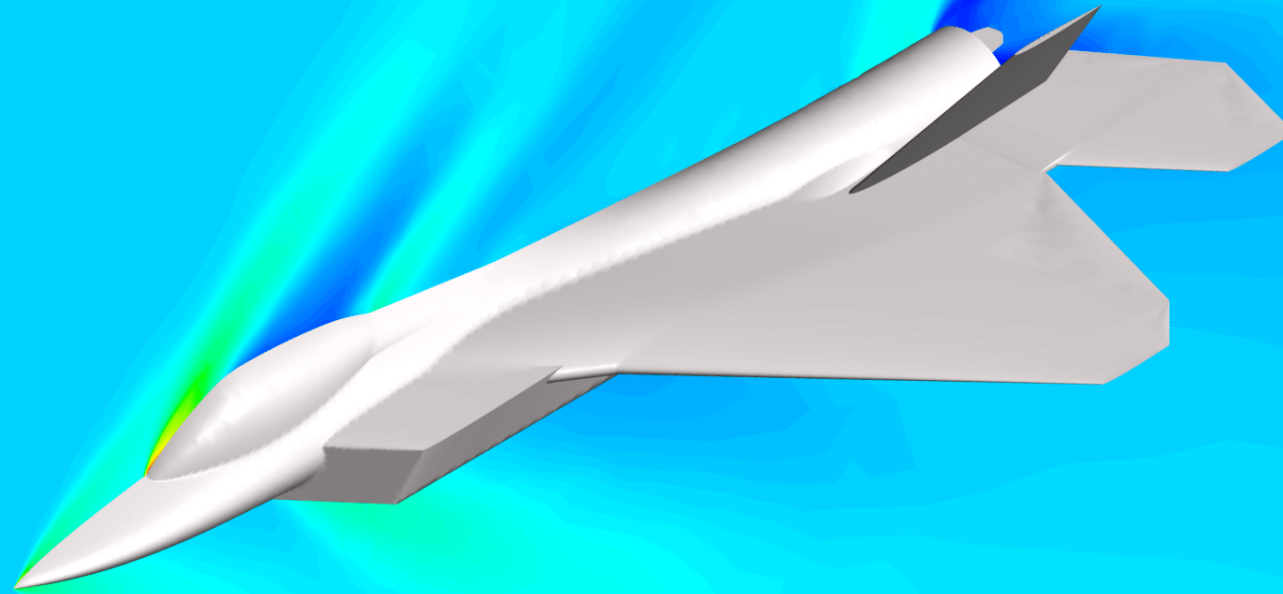
  
2      2      2

Number of steps  $\propto$  number of dimensions

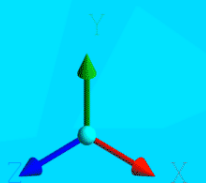
Pressure  
Contour 1



[Pa]



**What if  $f$  is very time-consuming or expensive to compute?**





## What if $n$ is very very big?



27 billion (Gemma 2, 2025)



671 billion (R1, 2025)



**ChatGPT**

Billions? Trillions?

Number of parameters  $n$

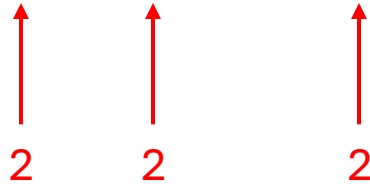
We need a better method of calculating gradients.

# Steps taken

$$\text{In 2D: } \nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

$$\text{In 3D: } \nabla f = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right)$$

$$\text{In nD: } \nabla f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$



Finite differences: Number of  $f$  calculations  $\propto$  number of dimensions

Automatic differentiation will give the gradient **in at most 4x the number of forward pass (function  $f$ ) operations<sup>1</sup>!**

[1] A. Griewank and A. Walther, Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation, 2nd ed., SIAM, Philadelphia, 2008

# Automatic/algorithmic differentiation

using the information *already present in your code* to calculate the gradient.

*Not a new idea:*

*1952 Master's Thesis from John F. Nolan, Boston University: Analytical Differentiation on a digital computer*

$$f(x_1, x_2, x_3) = (x_1 - x_2 + x_2x_3) * (\cos x_2x_3)$$

Find  $\frac{\partial f}{\partial x_1}$ ,  $\frac{\partial f}{\partial x_2}$ ,  $\frac{\partial f}{\partial x_3}$

$$f(x_1, x_2, x_3) = \overbrace{(x_1 - x_2 + x_2x_3)}^{z_1} * \overbrace{(\cos x_2x_3)}^{z_2}$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial z_1}{\partial x_1} z_2 + z_1 \frac{\partial z_2}{\partial x_1} \quad \text{By product rule}$$

$$z_1 = \underbrace{x_1 - x_2}_{y_1} + \underbrace{x_2x_3}_{y_2}$$



$$z_2 = \cos \underbrace{x_2x_3}_{y_2}$$

$$\begin{aligned} \frac{\partial z_1}{\partial x_1} &= \frac{\partial y_1}{\partial x_1} + \frac{\partial y_2}{\partial x_1} \\ &= 1 + 0 \end{aligned}$$

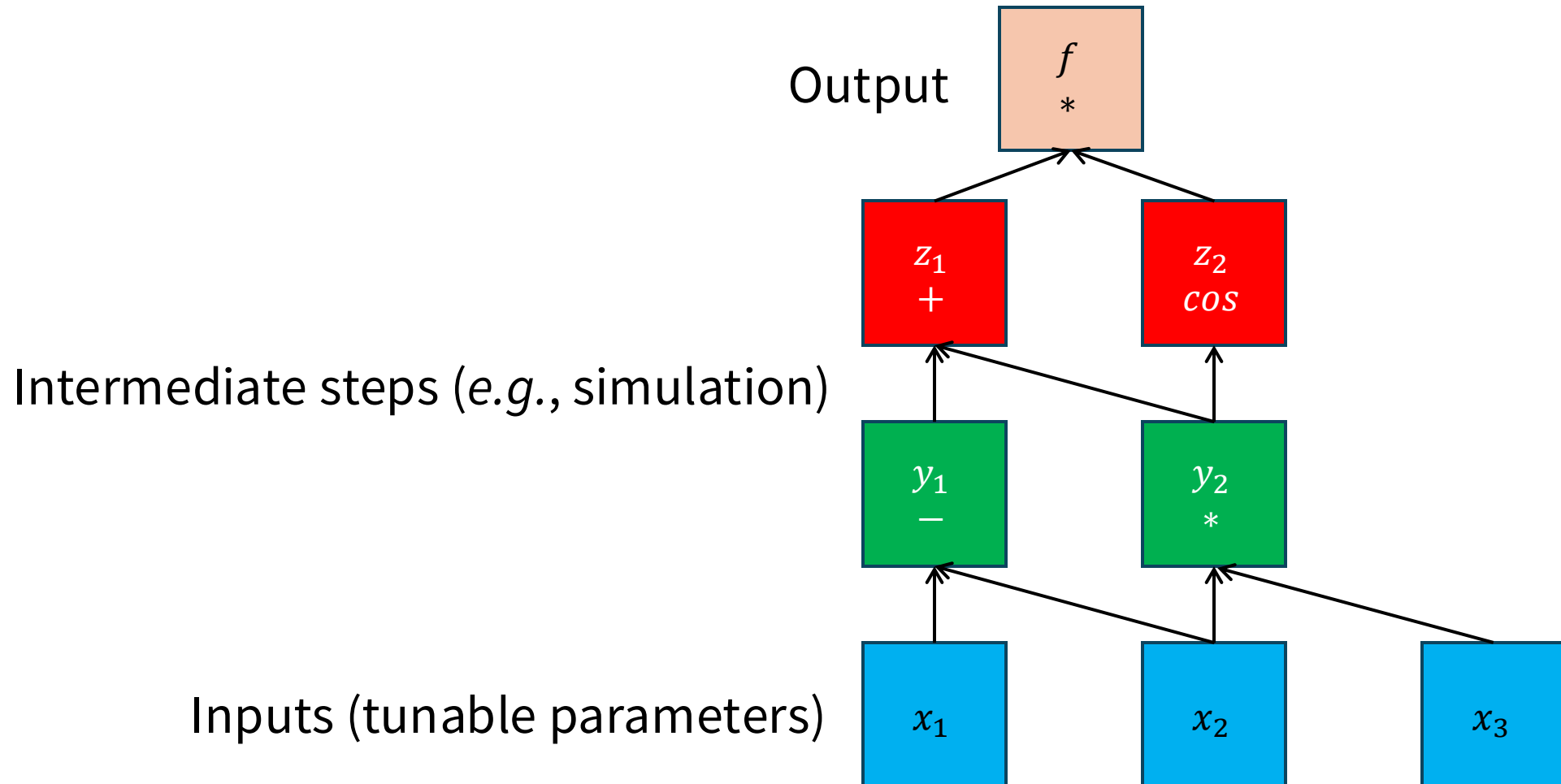
$$\begin{aligned} \frac{\partial z_2}{\partial x_1} &= -\sin y_2 * \frac{\partial y_2}{\partial x_1} \quad \text{By chain rule} \\ &= 0 \end{aligned}$$

$$\frac{\partial f}{\partial x_1} = \cos x_2x_3$$

# Equivalent computational tree

$$f(x_1, x_2, x_3) = \underbrace{(x_1 - x_2 + x_2 x_3)}_{z_1} * \underbrace{(\cos x_2 x_3)}_{z_2}$$

*Note: In the original image, the expression is annotated with green brackets for intermediate steps  $y_1$  and  $y_2$ , and red brackets for final intermediate results  $z_1$  and  $z_2$ .*





$$f(x_1, x_2, x_3) = \overbrace{(x_1 - x_2 + x_2x_3)}^{z_1} * \overbrace{(\cos x_2x_3)}^{z_2}$$

$$\frac{\partial f}{\partial x_1} = \frac{\partial z_1}{\partial x_1} z_2 + z_1 \frac{\partial z_2}{\partial x_1} \quad \text{By product rule}$$

$$z_1 = \underbrace{x_1 - x_2}_{y_1} + \underbrace{x_2x_3}_{y_2}$$



$$z_2 = \cos \underbrace{x_2x_3}_{y_2}$$

$$\begin{aligned} \frac{\partial z_1}{\partial x_1} &= \frac{\partial y_1}{\partial x_1} + \frac{\partial y_2}{\partial x_1} \\ &= 1 + 0 \end{aligned}$$

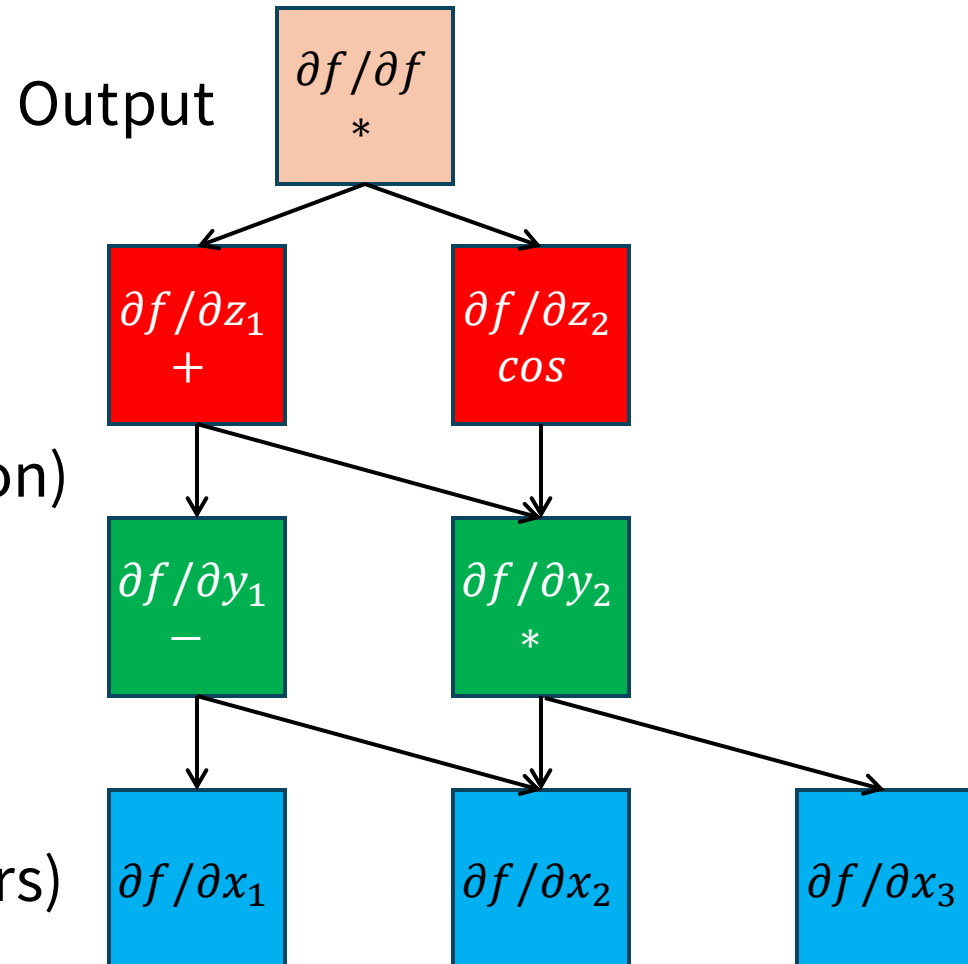
$$\begin{aligned} \frac{\partial z_2}{\partial x_1} &= -\sin y_2 * \frac{\partial y_2}{\partial x_1} \quad \text{By chain rule} \\ &= 0 \end{aligned}$$

$$\frac{\partial f}{\partial x_1} = \cos x_2x_3$$

# Adjoint computational tree

$$f(x_1, x_2, x_3) = \underbrace{(x_1 - x_2 + x_2 x_3)}_{z_1} * \underbrace{(\cos x_2 x_3)}_{z_2}$$

$y_1$                        $y_2$                        $y_2$



# Key insight for autodiff

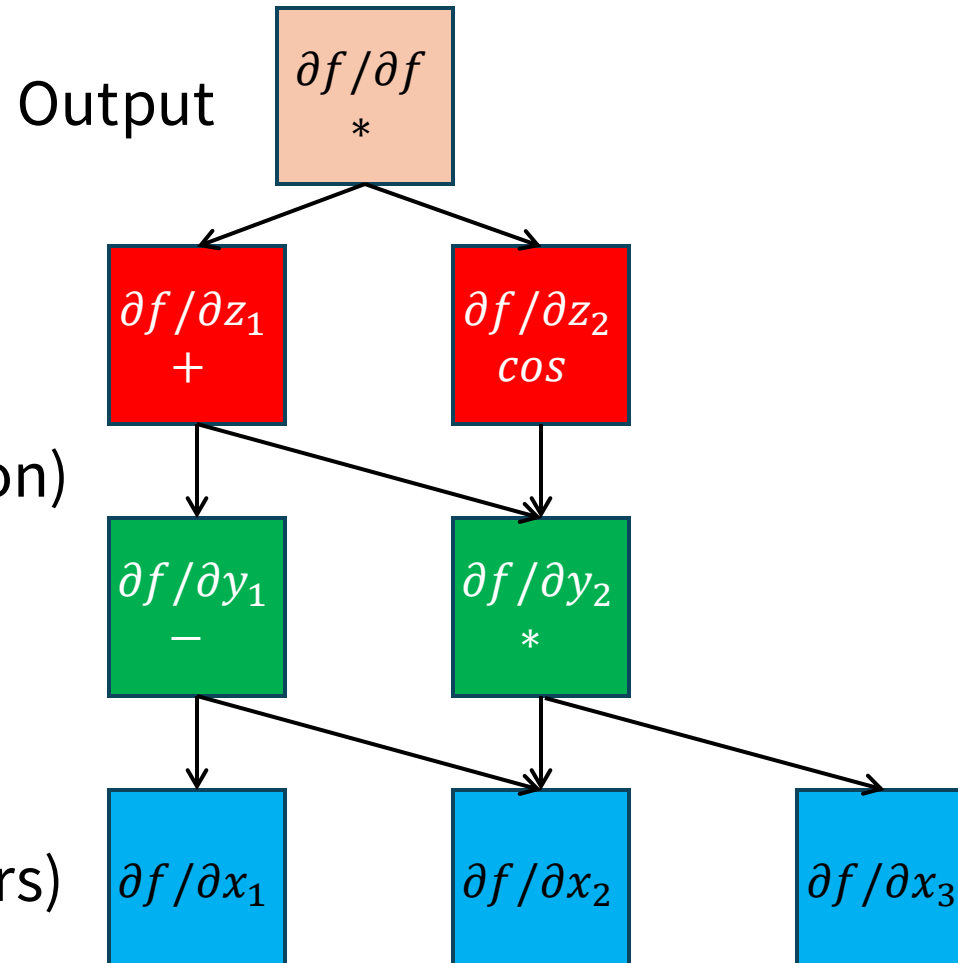
The sensitivity  $\frac{\partial f}{\partial j}$  of every node  $j$  can be computed by tracing the computational tree backwards.



*How does the output change by wiggling this value?*

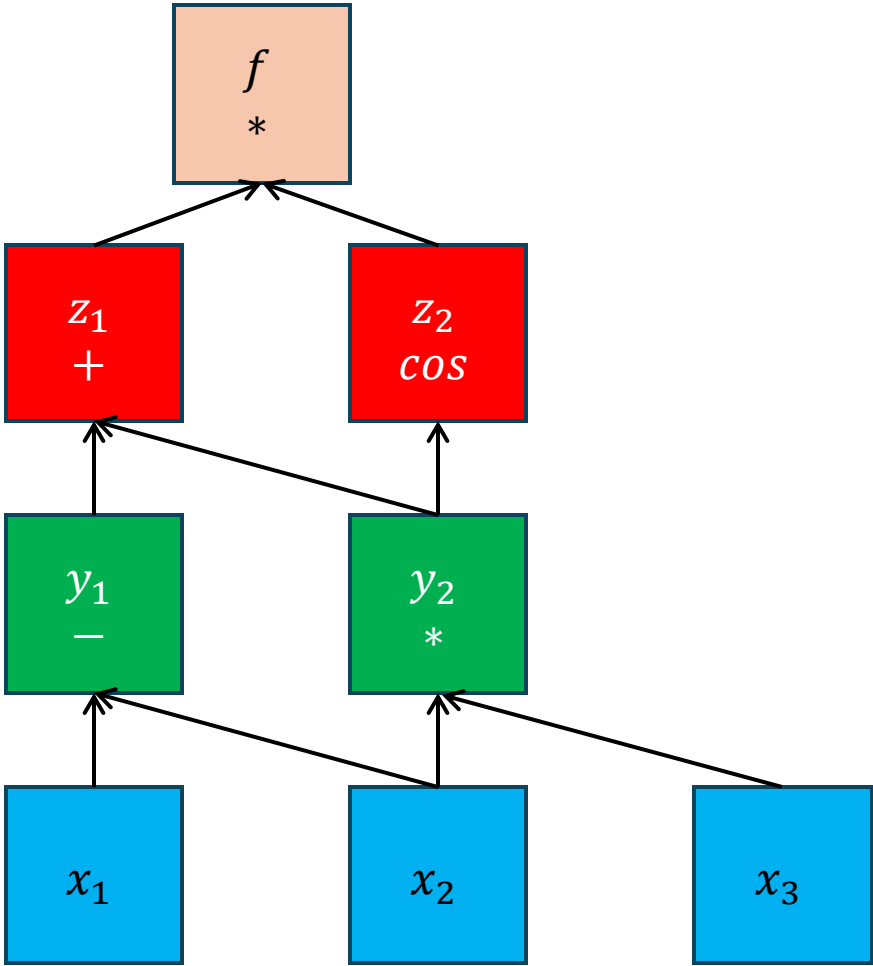
Intermediate steps (e.g., simulation)

Inputs (tunable parameters)

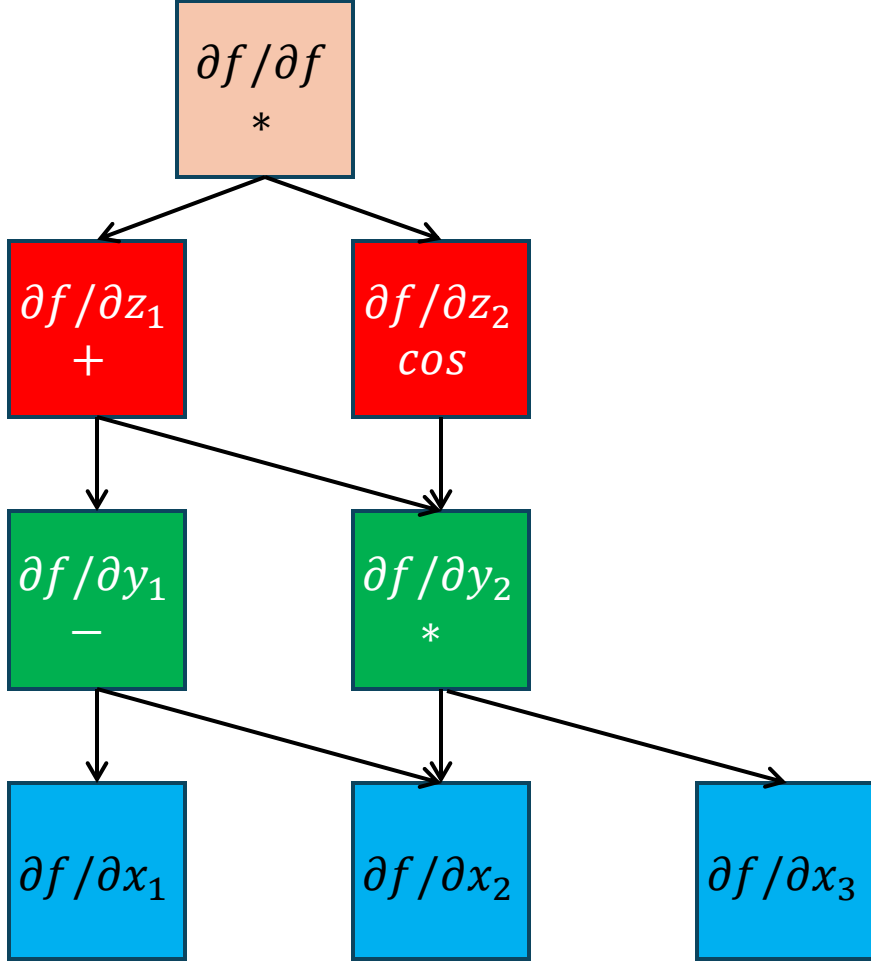


# Autodiff protocol

1. Run forward pass. Store computational tree and intermediate values.



2. Run reverse pass. Use stored tree to compute sensitivities.



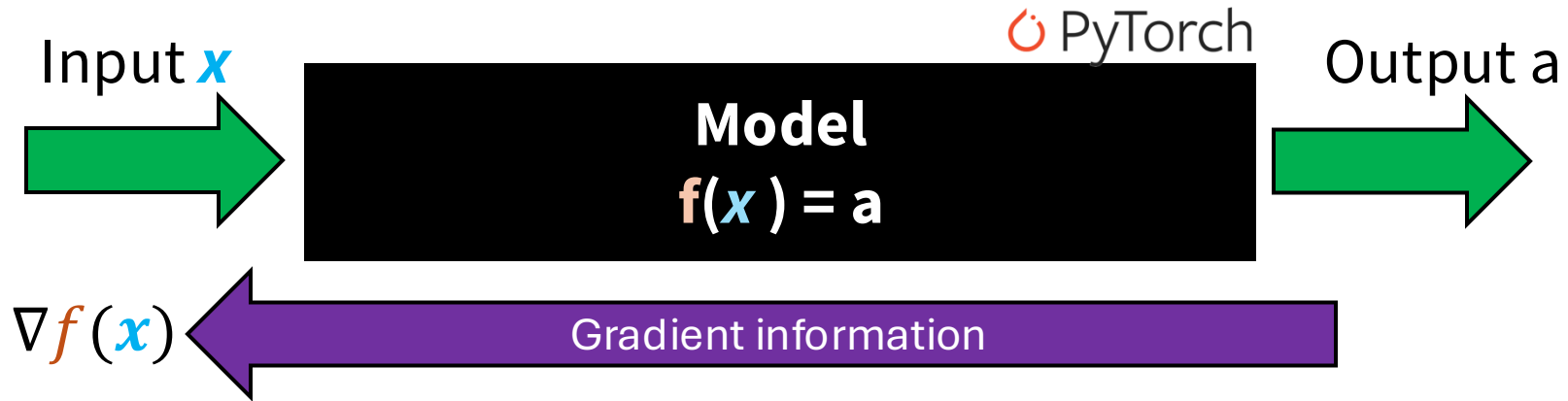
Do you need to know about computational trees to use Autodiff?

No (: Free frameworks build it for you automatically!

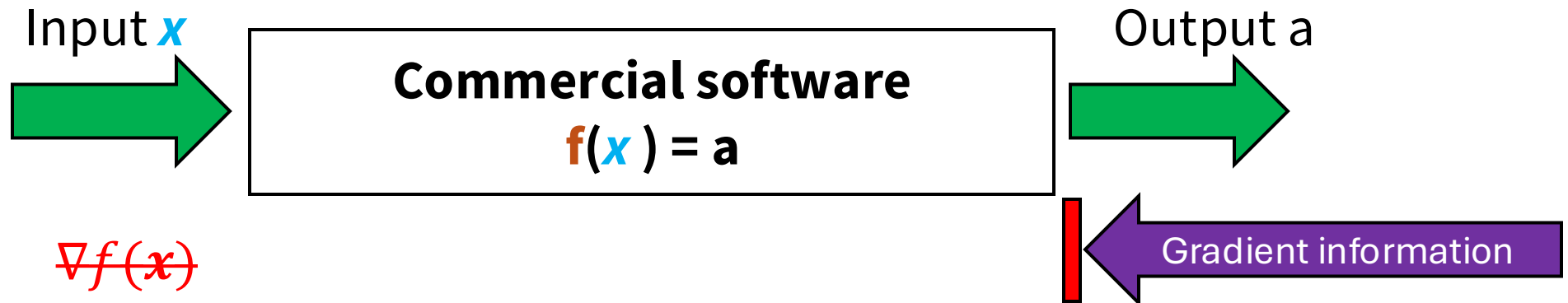


# A small catch

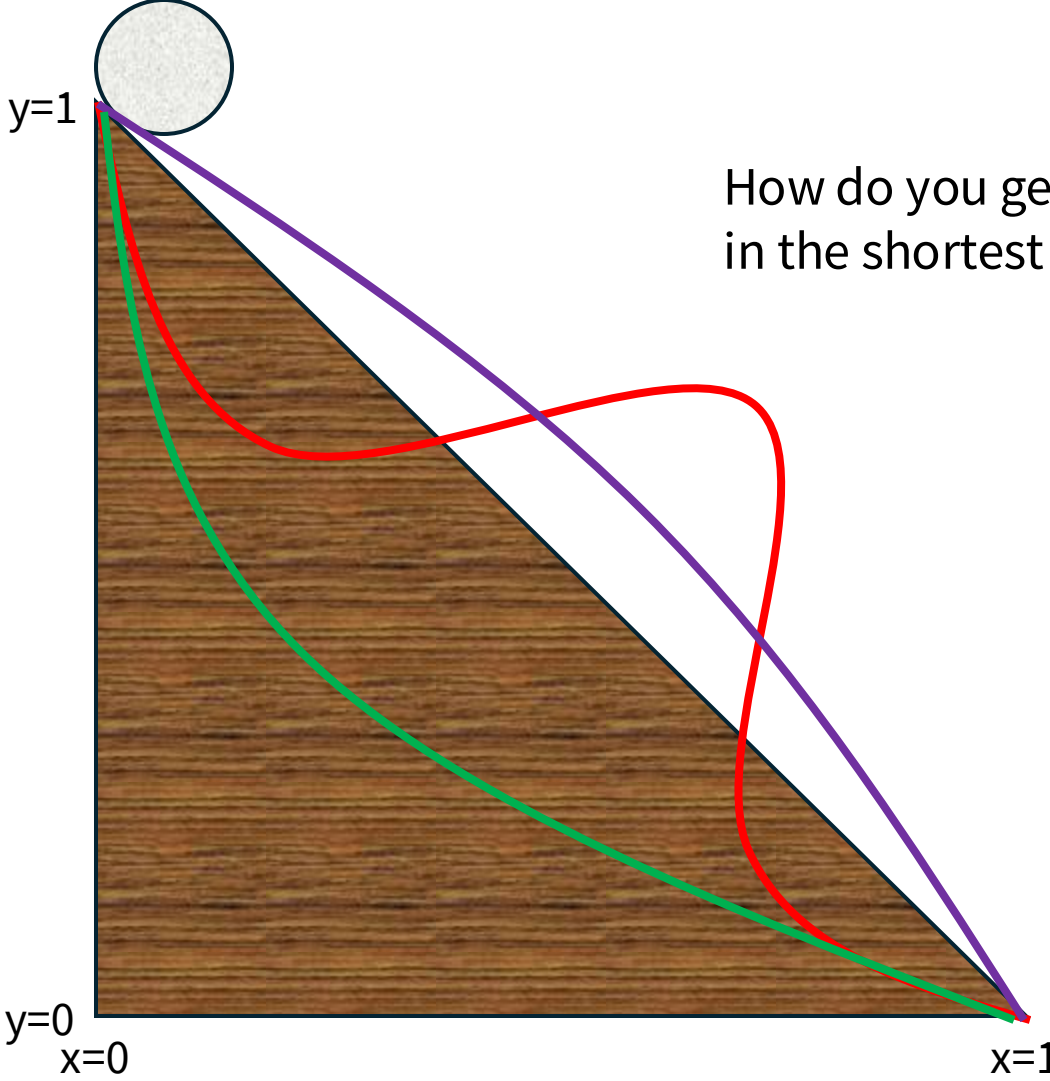
All calculations need to be performed on a *differentiable platform* for the gradients to be backpropagated.



Software lacking the source code are *not differentiable*.



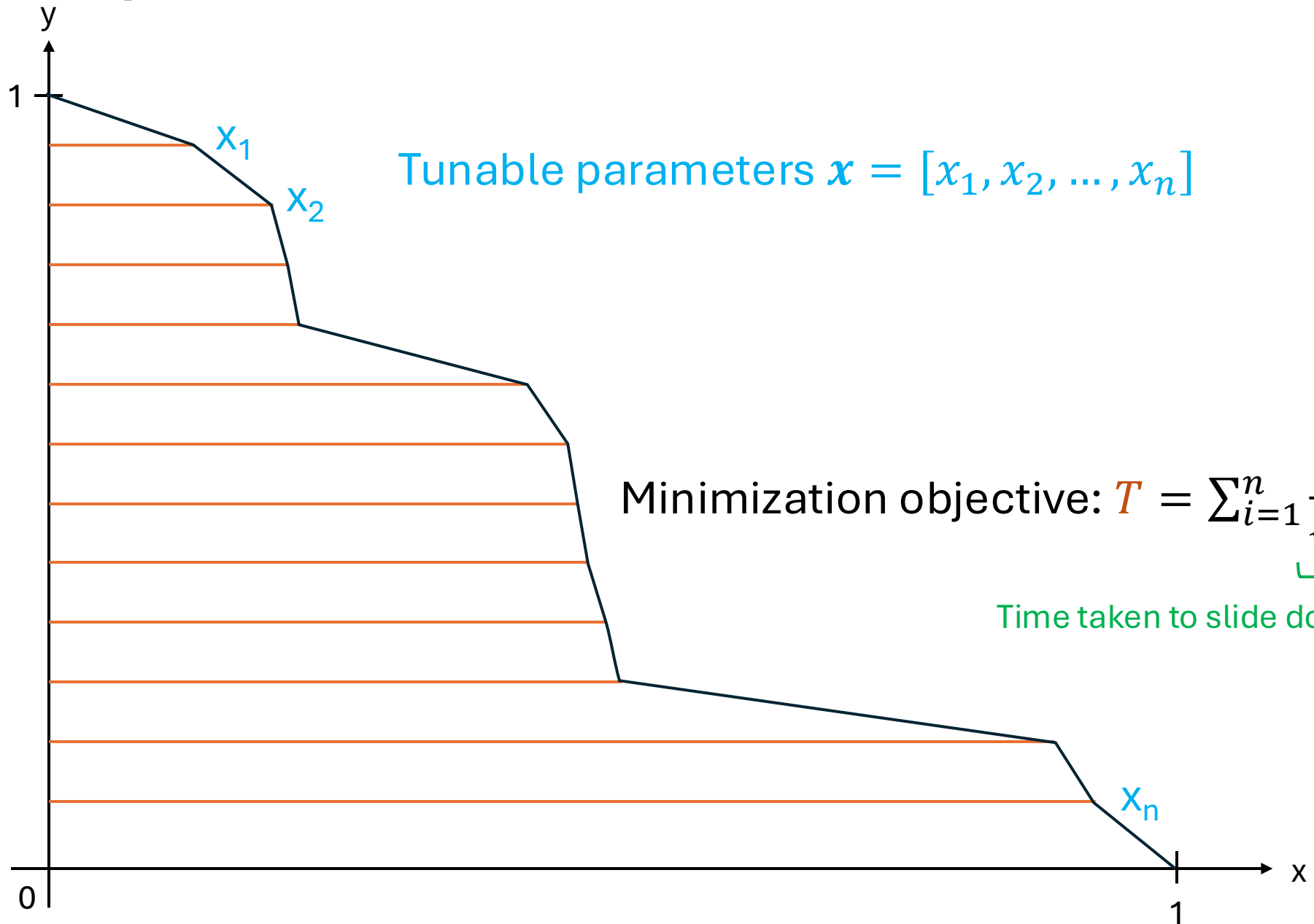
# Let's solve a problem



How do you get a particle to slide to the end in the shortest possible time?



# Optimization setup



# Recap

- The time-intensive step in optimization is often gradient calculation.
- Automatic differentiation enables efficient high-dimensional gradient calculation for any physical or mathematical system.
- To use automatic differentiation, all calculations must be on a *differentiable* platform.

Download these slides and demo code here:



<https://danlimsw.com/coursenotes/>