# INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3
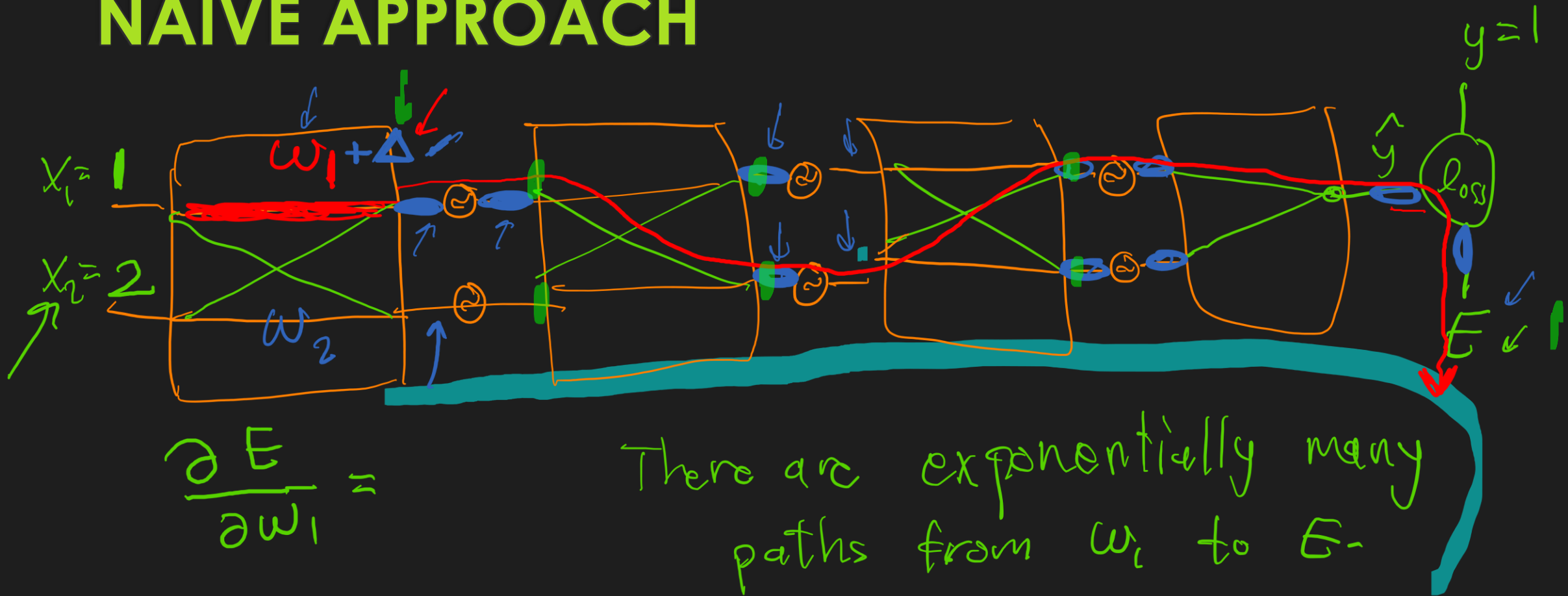
## LECTURE 23

HASSAN ASHTIANI

# COMPUTING THE GRADIENT EFFICIENTLY

- $\nabla_w\big(\mathrm{L}(w,b)\big) = \sum_i \nabla_w\left(l\big(f_{w,b}(x^i), y^i\big)\right)$

  - How to calculate $\nabla_w\left(l\big(f_{w,b}(x^i), y^i\big)\right)$?

- This can be computationally expensive

# NAÏVE APPROACH



$x_1 = 1$

$x_2 = 2$

$\omega_1 + \Delta$

$\omega_2$

$y = 1$

$\hat{y}$

loss

$E$

$$\frac{\partial E}{\partial \omega_1} =$$

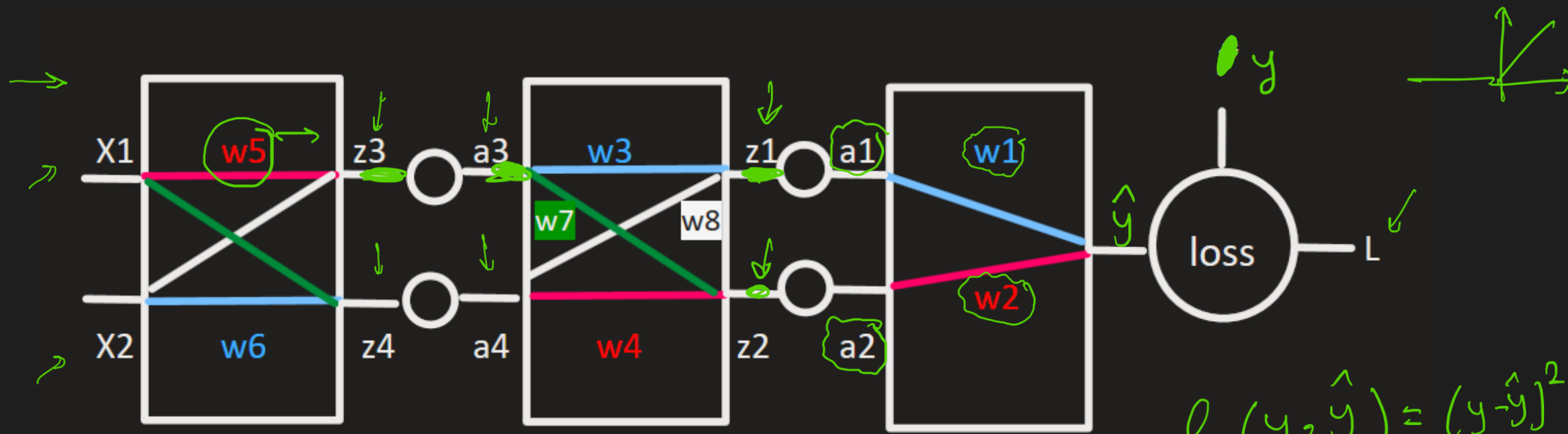There are exponentially many paths from $\omega_1$ to $E$.

* the computations involved in computing $\frac{\partial E}{\partial \omega_1}$ and $\frac{\partial E}{\omega_2}$ are redundant

# COMPUTING THE GRADIENT?

- Naïve approach:
  - Computationally expensive for deep models
  - Some of the computations are repetitive
- What to do?
  - A kind of "dynamic programming"
  - Back propagation

# BACK-PROPAGATION

- Use chain rule (for vector-valued functions)

- Linear time in terms of the number of weights!

- Forward phase

  - Compute the input/outputs of activation functions

- Backward phase

  - Compute the gradients, layer-by-layer

$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

ReLU activation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = 2(\hat{y} - y) \cdot a_1$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} = 2(\hat{y} - y) \cdot w_1 \cdot \mathbb{1}\{z_1 > 0\}$$

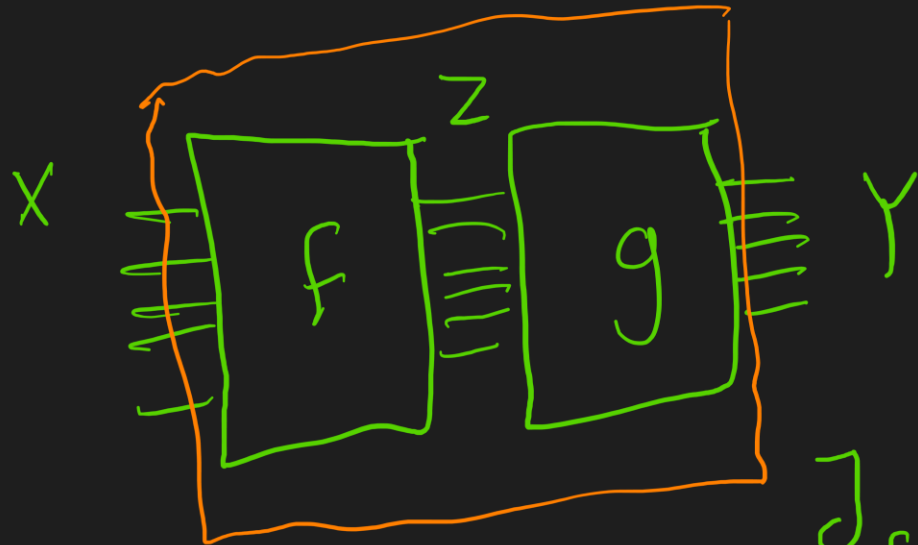$$\frac{\partial L}{\partial w_5} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial w_5} =$$

$$\left( \frac{\partial L}{z_1} \cdot \frac{\partial z_1}{\partial a_3} + \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_3} \right) \cdot \frac{\partial a_3}{\partial z_3} \cdot x_1$$

$$=$$

How to do these using matrix/vector calculus?

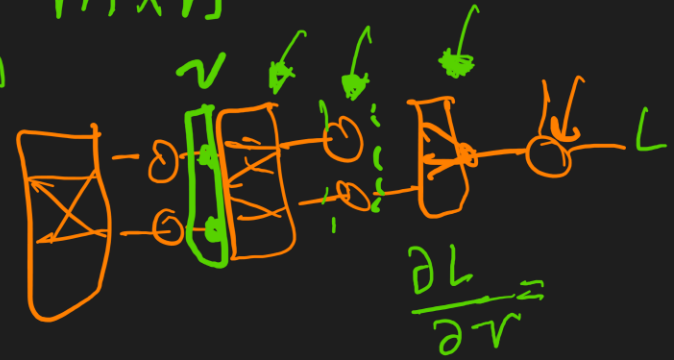# MULTI-VARIATE CHAIN RULE
# AND A MODULAR APPROACH

$A, B, C, D, E$

$f: \mathbb{R}^n \longrightarrow \mathbb{R}^m$

$g: \mathbb{R}^m \longrightarrow \mathbb{R}^d$

$$J_f = \begin{bmatrix} \dfrac{\partial f_1}{\partial x_1} & \cdots & \dfrac{\partial f_1}{x_n} \\ \vdots & & \\ \dfrac{\partial f_m}{\partial x_1} & & \dfrac{\partial f_m}{\partial x_n} \end{bmatrix}_{m \times n}$$

$$\nabla_x f \circ g(x) = [\nabla g]_{d \times m} \cdot [\nabla f]_{m \times n}$$

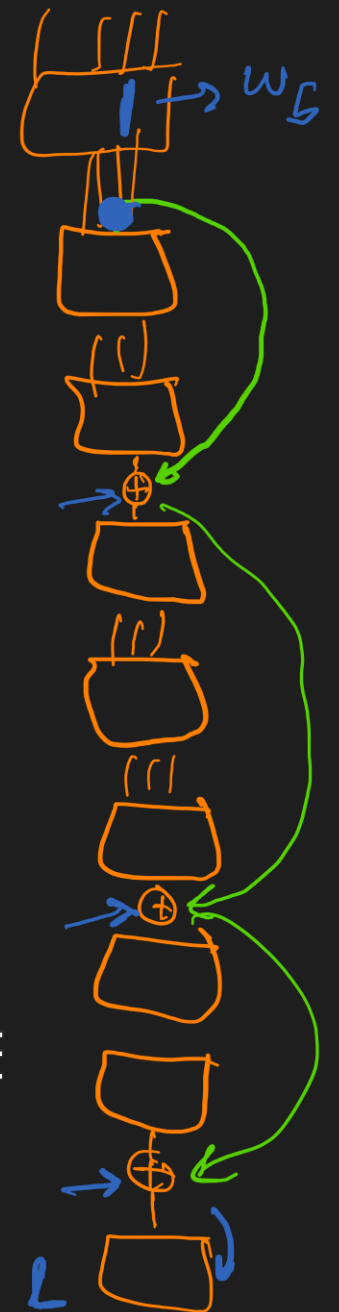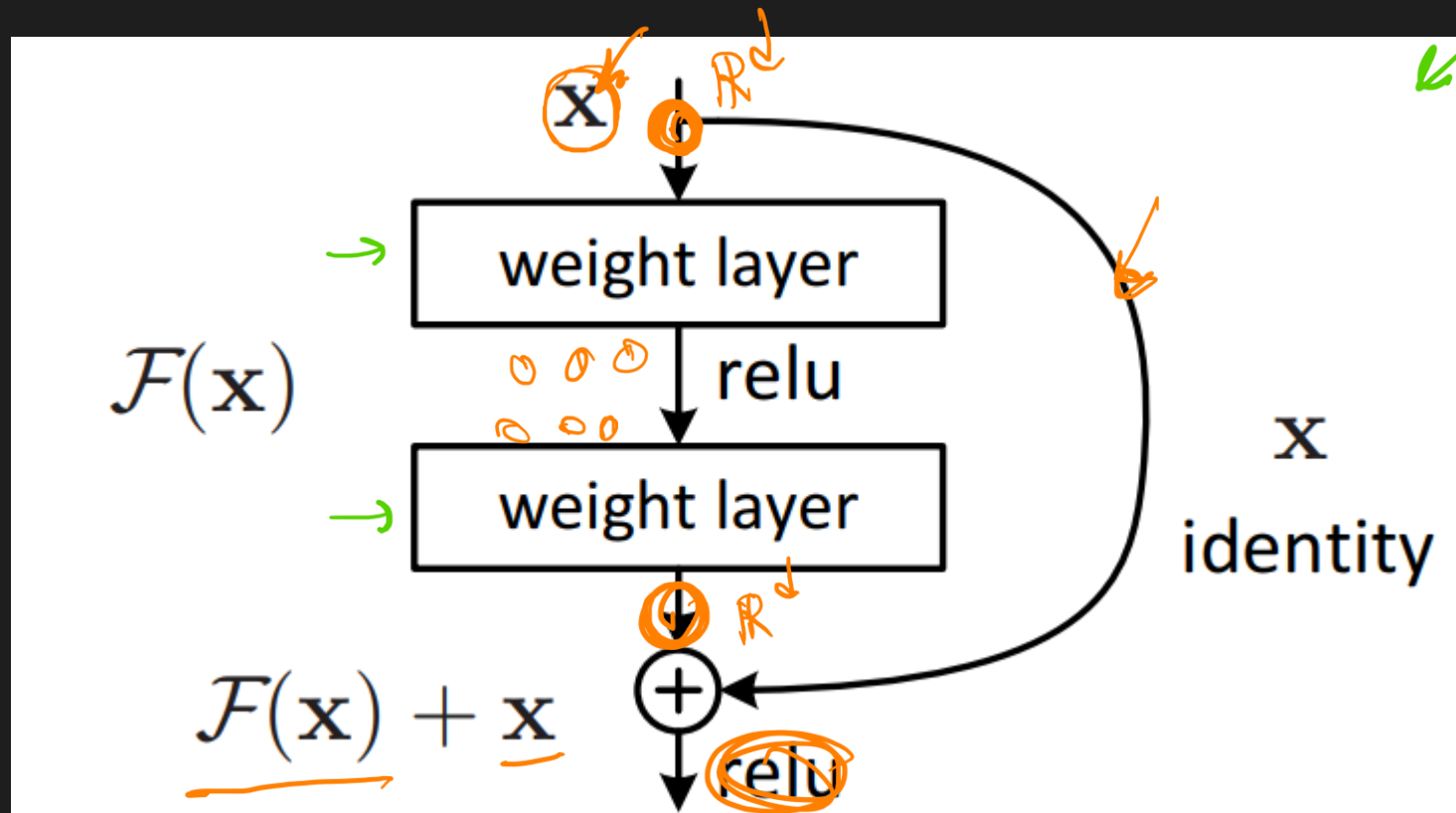if $f = Ax$ then $\nabla_f = A$

# VANISHING GRADIENT

- FOR DEEPER NETWORKS, THE PARTIAL DERIVATIVES
  - FOR THE OUTPUT LAYER VS THE INPUT LAYER?
- FOR SIGMOID ACTIVATION FUNCTIONS
  - $\sigma'(x) \in [0,1]$
  - FOR "SATURATED" NEURONS $\sigma'(x) \approx 0$
  - GRADIENT "DOES NOT REACH" THE FIRST LAYERS
- WHAT CAN WE DO?

# VANISHING GRADIENT

- Some activation functions are better
  - Leaky ReLU (also ReLU/max-out)
  - But gradient still can vanish after a couple layers
- Better initialization
- A possible workaround
  - Use "shortcuts" for the gradient to flow
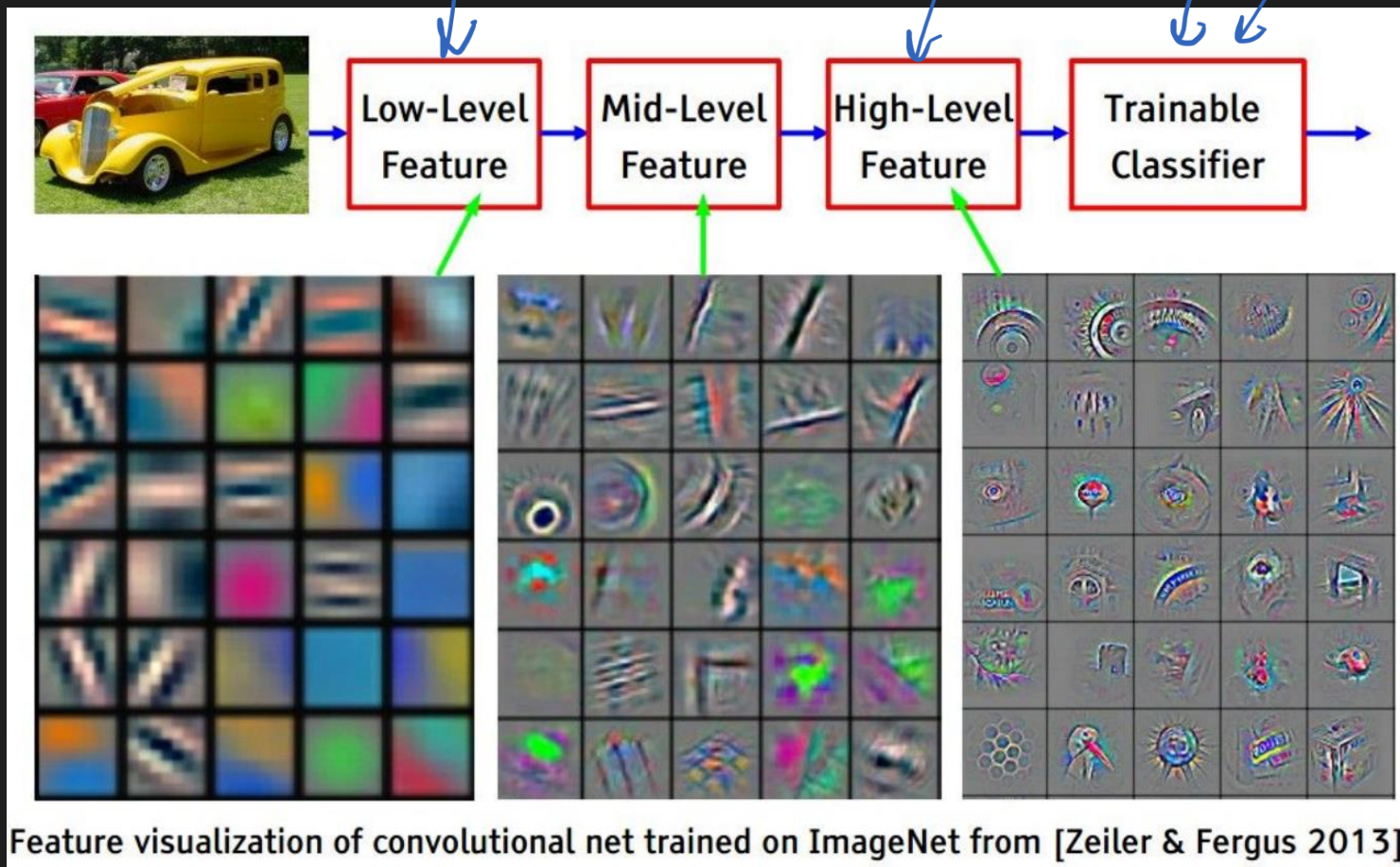- ResNets can have 100s of layers!

# DEEP RESIDUAL NETWORKS



- THE INPUT CAN DIRECTLY GO TO DEEPER NEURONS, SO THE GRADIENT CAN FLOW

# UNIVERSAL APPROXIMATION THEOREM

- FEED-FORWARD NETWORKS WITH SIGMOID ACTIVATION FUNCTIONS CAN APPROXIMATE ANY BOUNDED CONTINUOUS FUNCTION UP TO DESIRABLE ACCURACY
    - **ONLY A SINGLE HIDDEN LAYER IS NEEDED!**
    - GEORGE CYBENKO, 1989
    - ALSO HOLDS FOR OTHER USUAL ACTIVATION FUNCTIONS
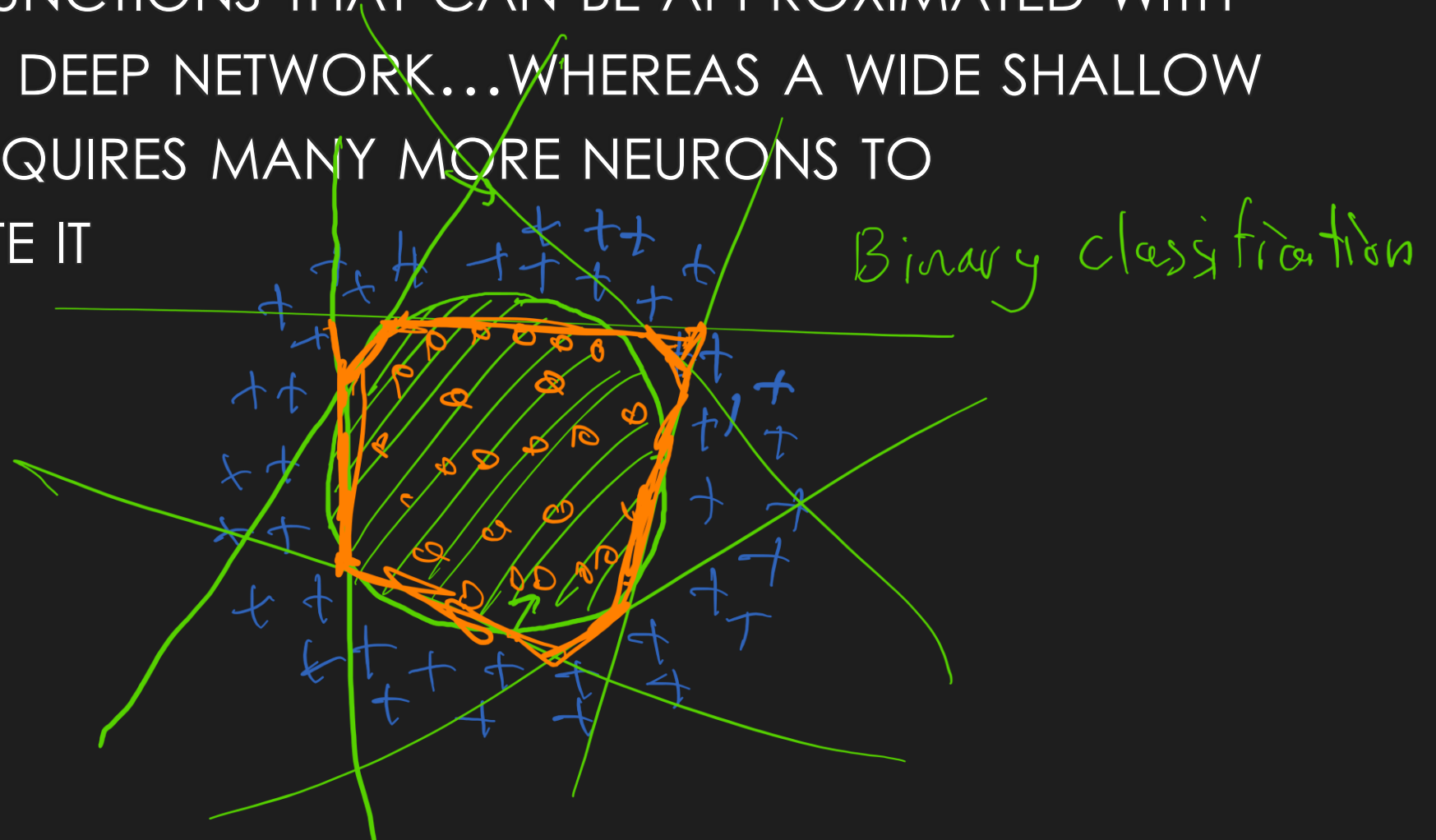- IS THERE A POINT IN HAVING MORE THAN ONE HIDDEN LAYER?

# DEEP VS SHALLOW NETWORKS
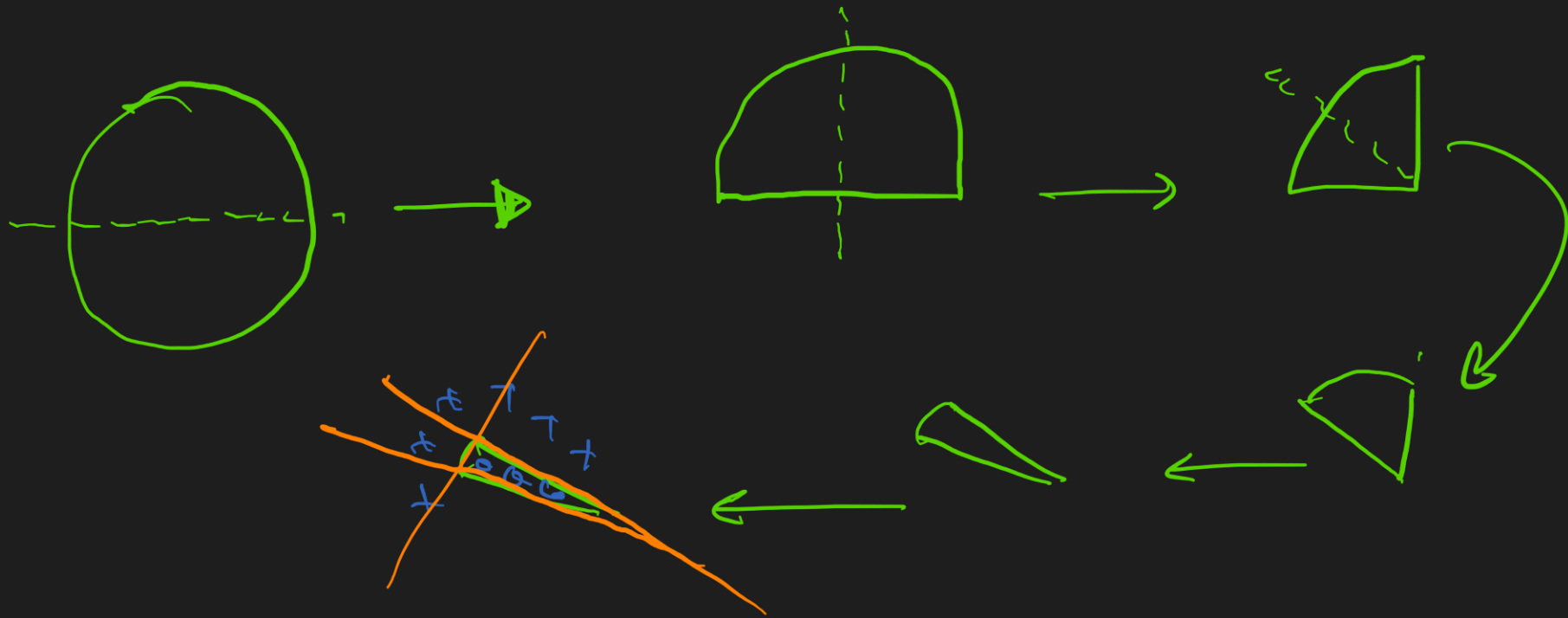
- Low-level to high-level computations/detections



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# DEEP VS SHALLOW NETWORKS

- THERE ARE FUNCTIONS THAT CAN BE APPROXIMATED WITH A SMALL BUT DEEP NETWORK…WHEREAS A WIDE SHALLOW NETWORK REQUIRES MANY MORE NEURONS TO APPROXIMATE IT

Binary classification

$$error \approx \frac{1}{\# neurons}$$

$$\text{error} \simeq \frac{1}{(\#\text{neurons})^{\text{layers}}}$$

# REGULARIZING NNS

- Good neural networks are often over-parametrized!
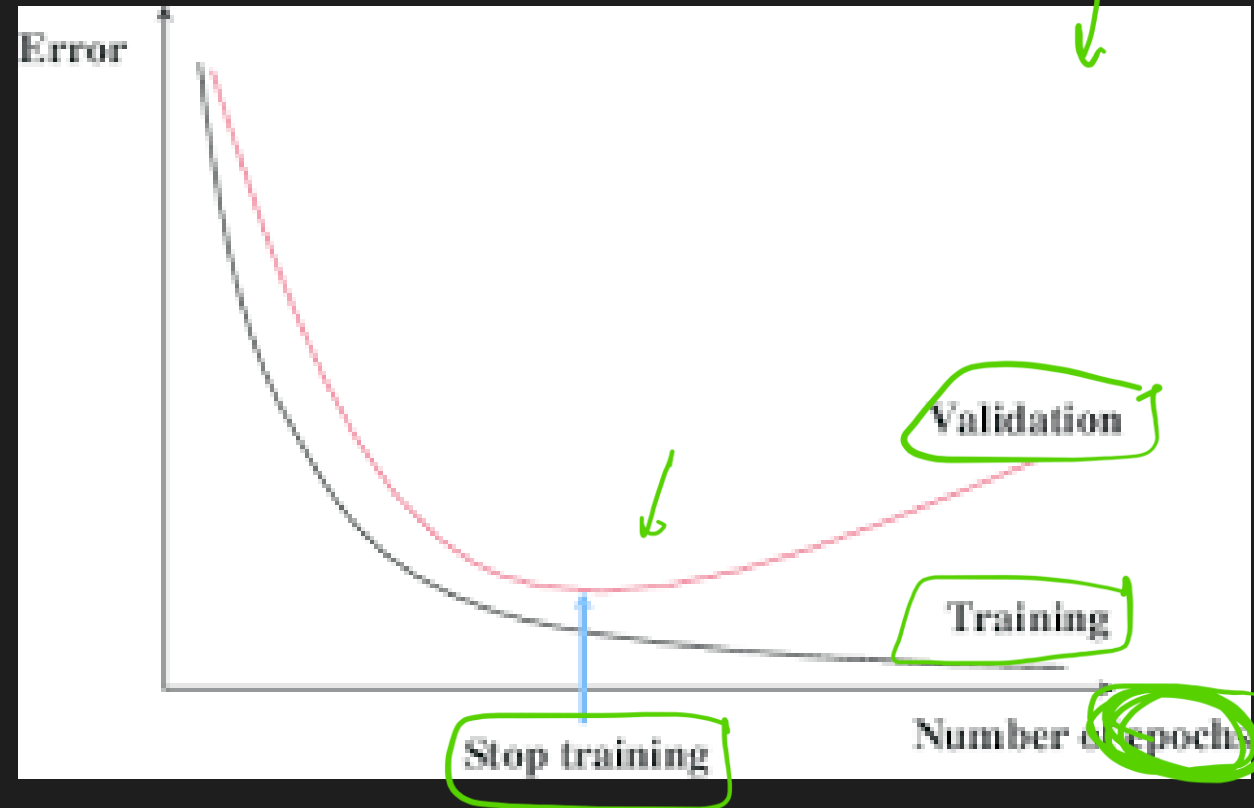  - Prone to over-fitting
1. Adding regularization terms to the objective function

  - $E(w) + \|w\|^2$
2. Early stopping
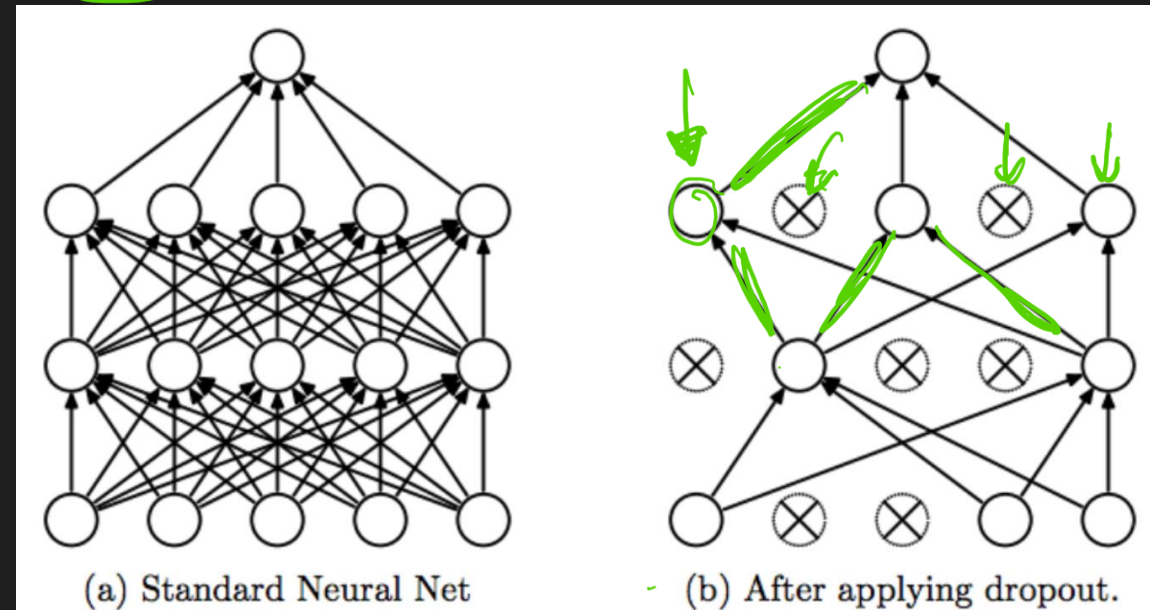3. Adding noise
4. Structural regularization
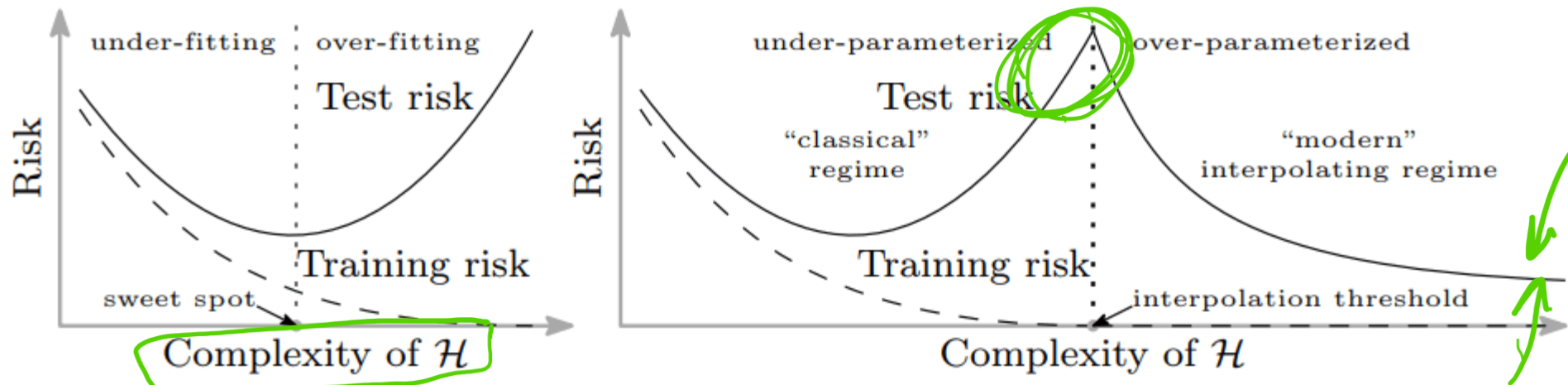
# REGULARIZING NNS WITH DROPOUT

TRAINING

- FOR EACH ITERATION OF STOCHASTIC GRADIENT DESCENT AND FOR EACH TRAINING DATA POINT DO:

  - DROP EACH NODE WITH PROBABILITY $p$

TESTING

- DON'T DROP THE NODES, BUT FOR ALL NODES, BUT MULTIPLY THE VALUE OF ACTIVATIONS BY $(1 - p)$



(a) Standard Neural Net      (b) After applying dropout.

# DOUBLE DESCENT?



(a) U-shaped "bias-variance" risk curve

(b) "double descent" risk curve