

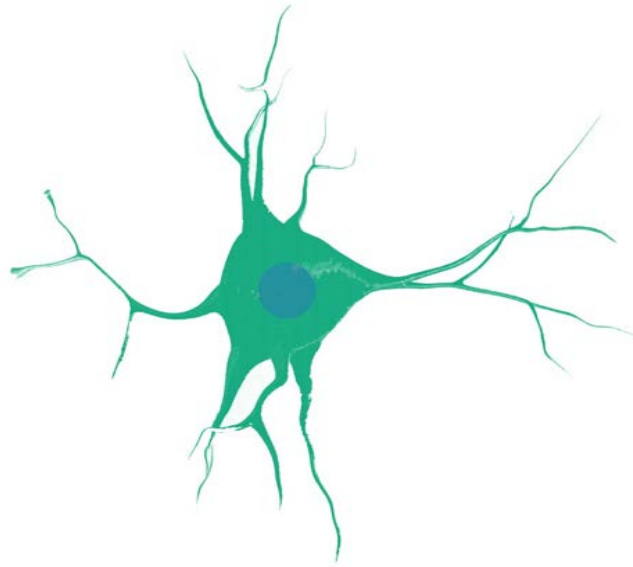
# Neural Networks

Cynthia Rudin

Duke Machine Learning

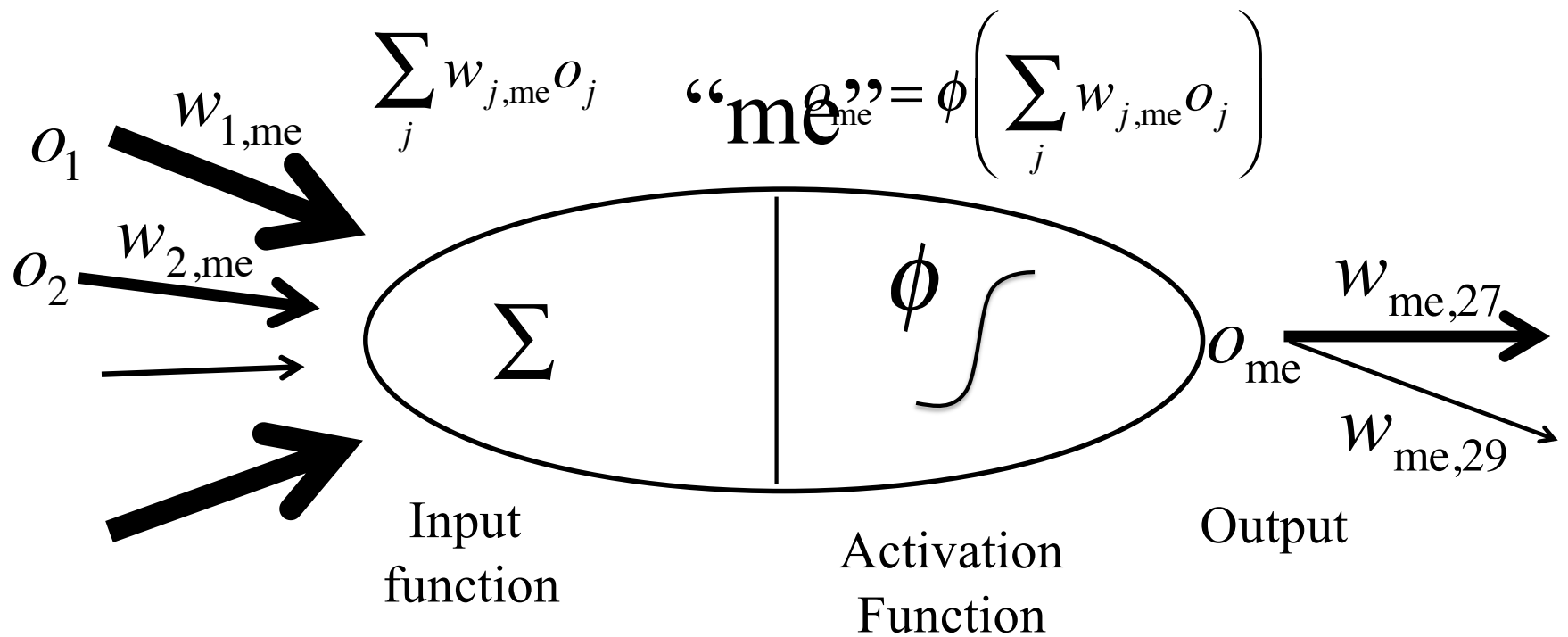
# Neurons

- $10^{11}$  neurons in a brain,  $10^{14}$  synapses (connections).
- Signals are electrical potential spikes that travel through the network.

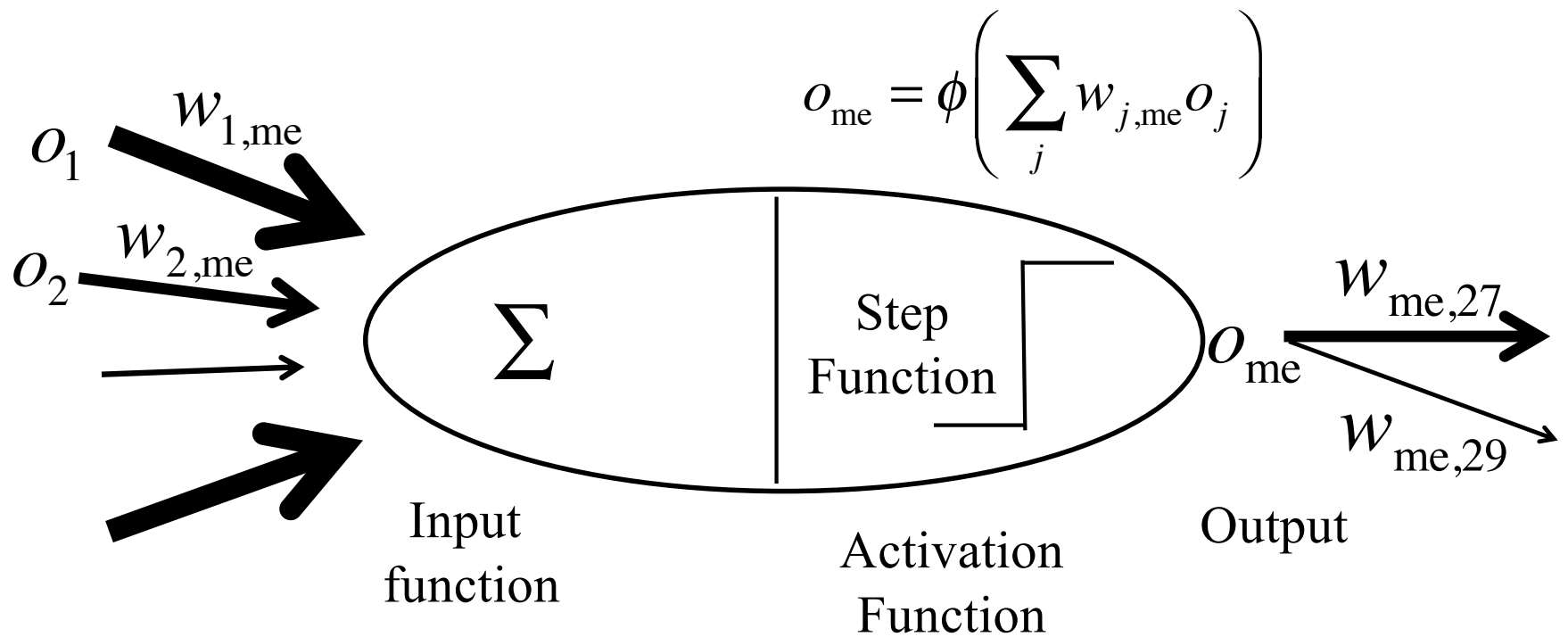


(Credit: Adapted from Russell and Norvig)

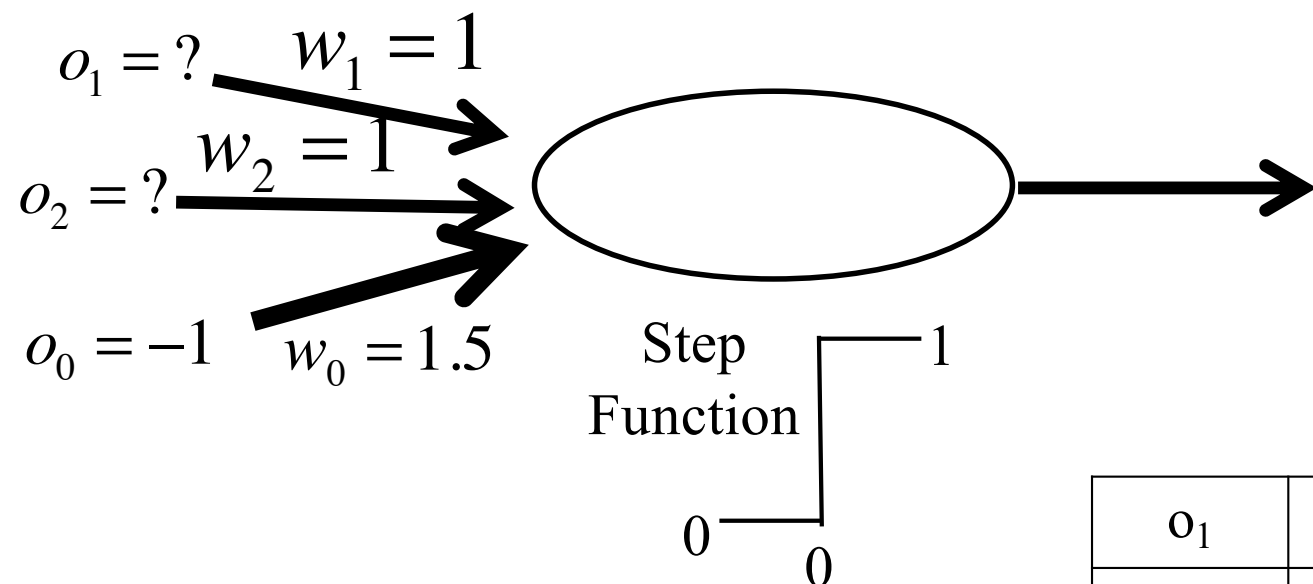
# McCulloch-Pitts “Neuron”



# McCulloch-Pitts “Neuron”

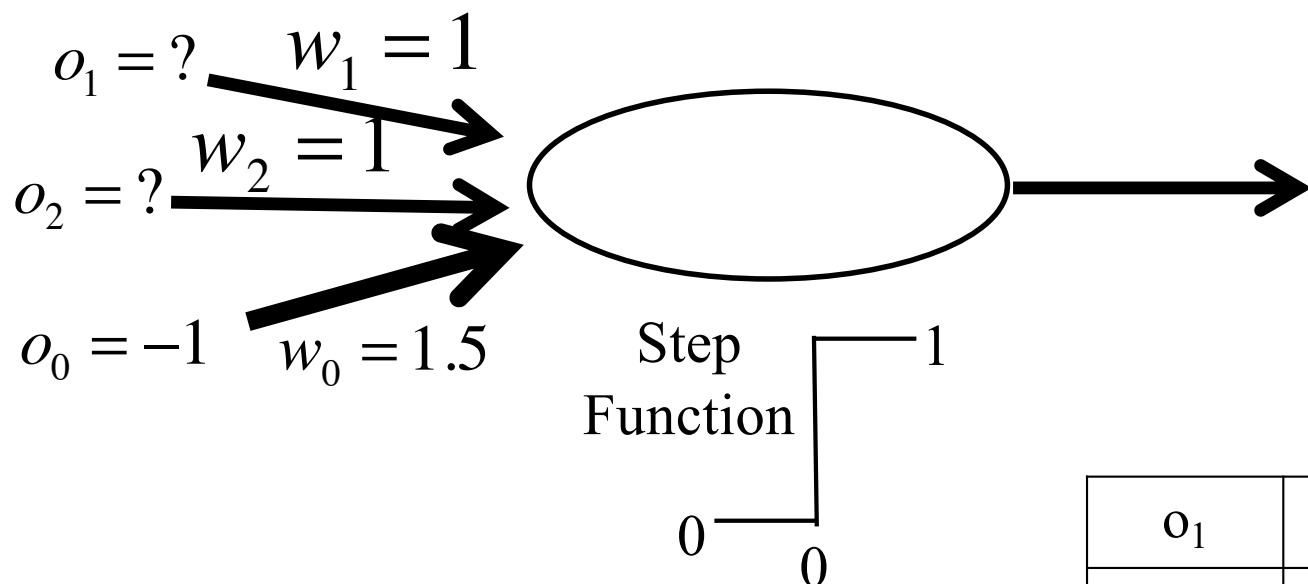


# McCulloch-Pitts “Neuron”



$o_1$	$o_2$	output

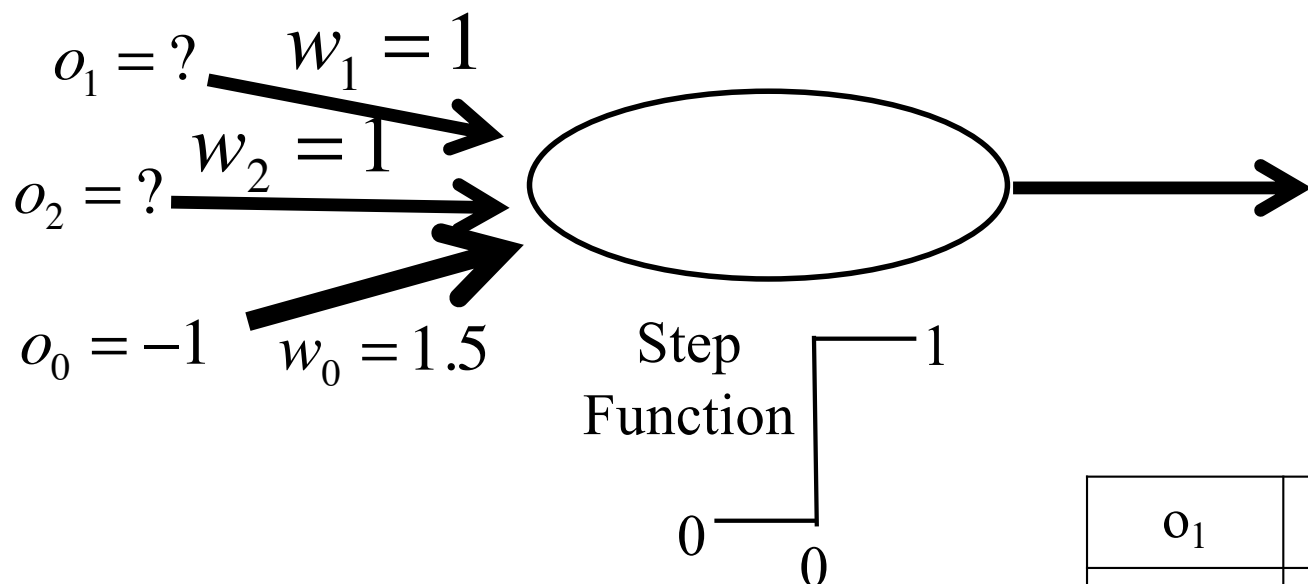
# McCulloch-Pitts “Neuron”



$$0 + 0 - 1.5 = -1.5$$

$o_1$	$o_2$	output
0	0	

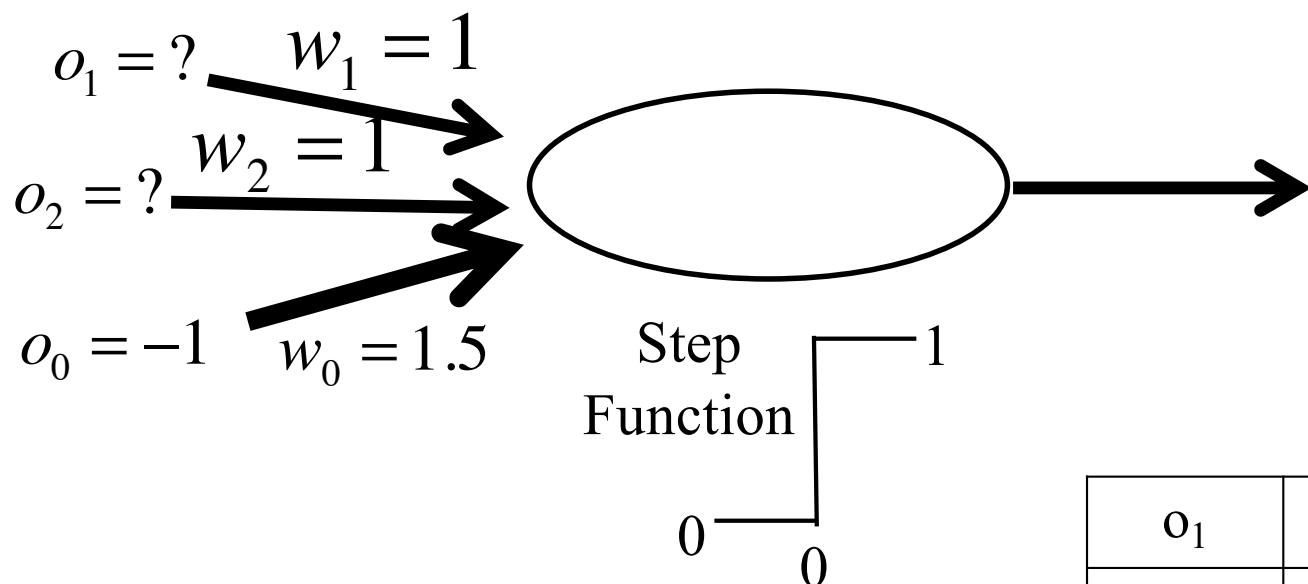
# McCulloch-Pitts “Neuron”



$$0 + 0 - 1.5 = -1.5$$

$o_1$	$o_2$	output
0	0	0

# McCulloch-Pitts “Neuron”

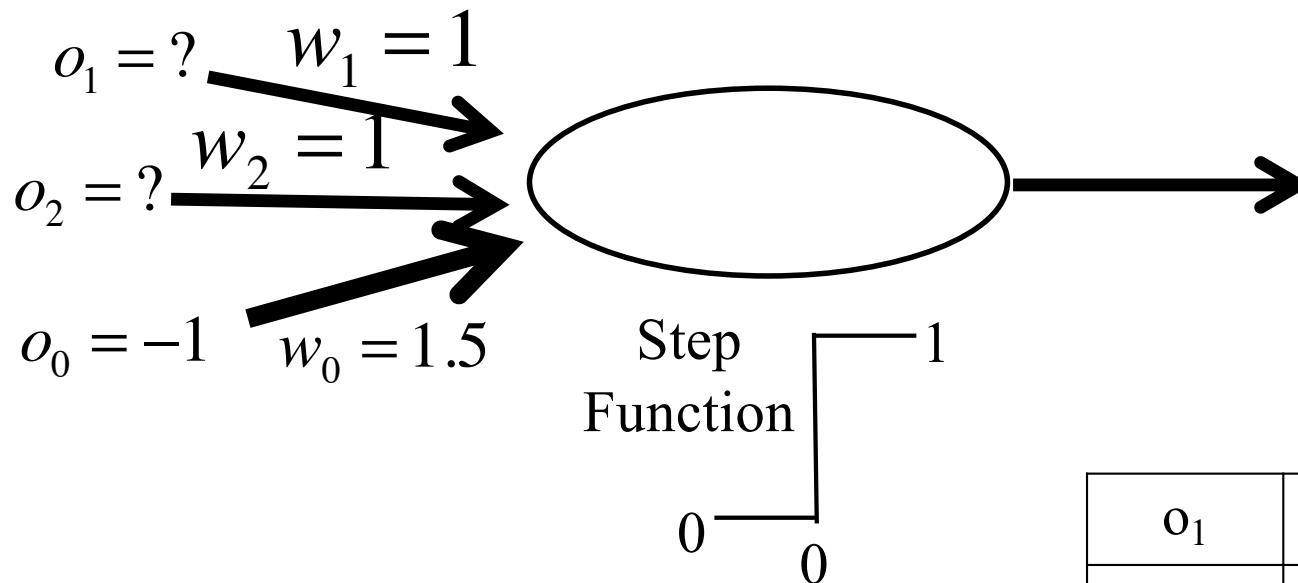


$$1 + 0 - 1.5 = -0.5$$

$o_1$	$o_2$	output
0	0	0
1	0	0



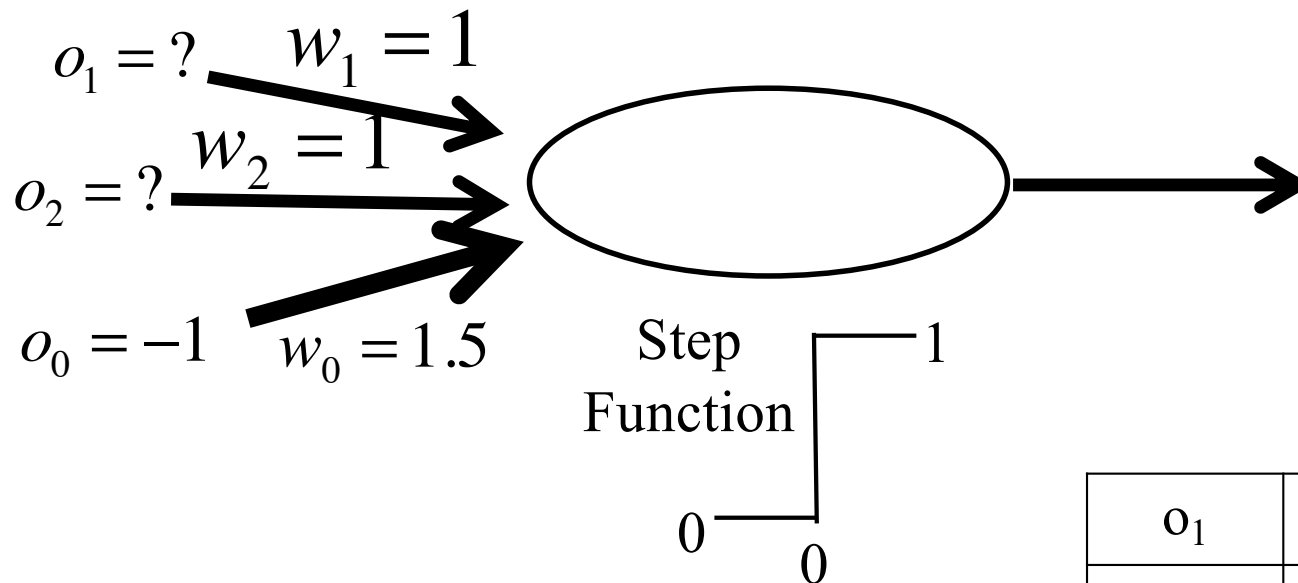
# McCulloch-Pitts “Neuron”



$$0 + 1 - 1.5 = -0.5$$

$o_1$	$o_2$	output
0	0	0
1	0	0
0	1	0

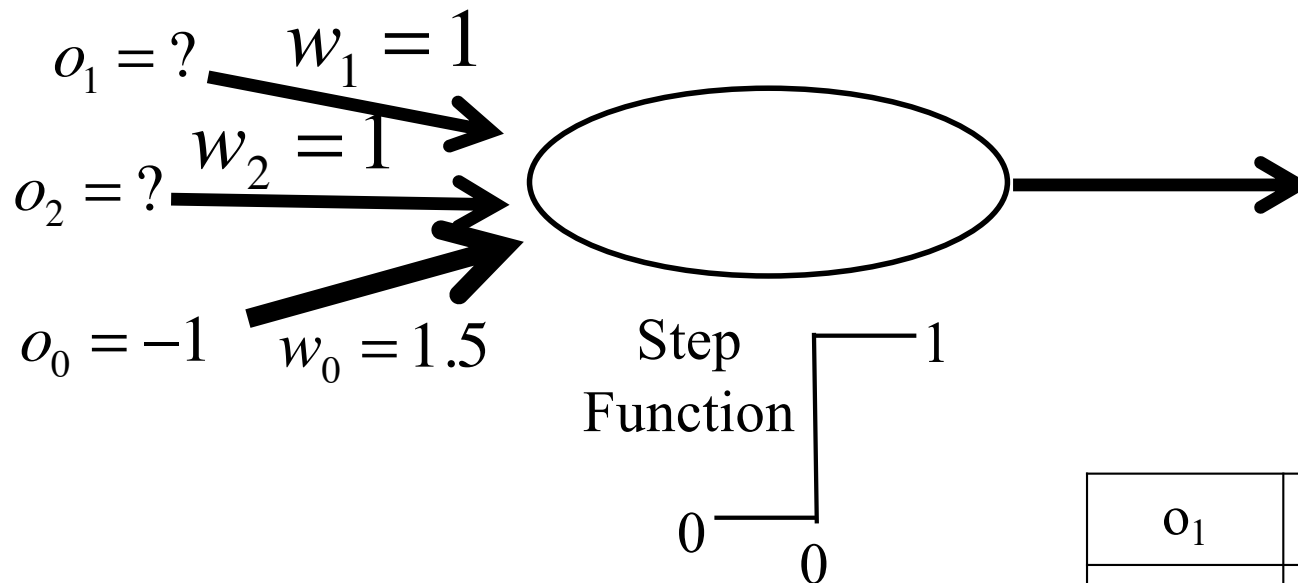
# McCulloch-Pitts “Neuron”



$$1 + 1 - 1.5 = 0.5$$

$o_1$	$o_2$	output
0	0	0
1	0	0
0	1	0
1	1	

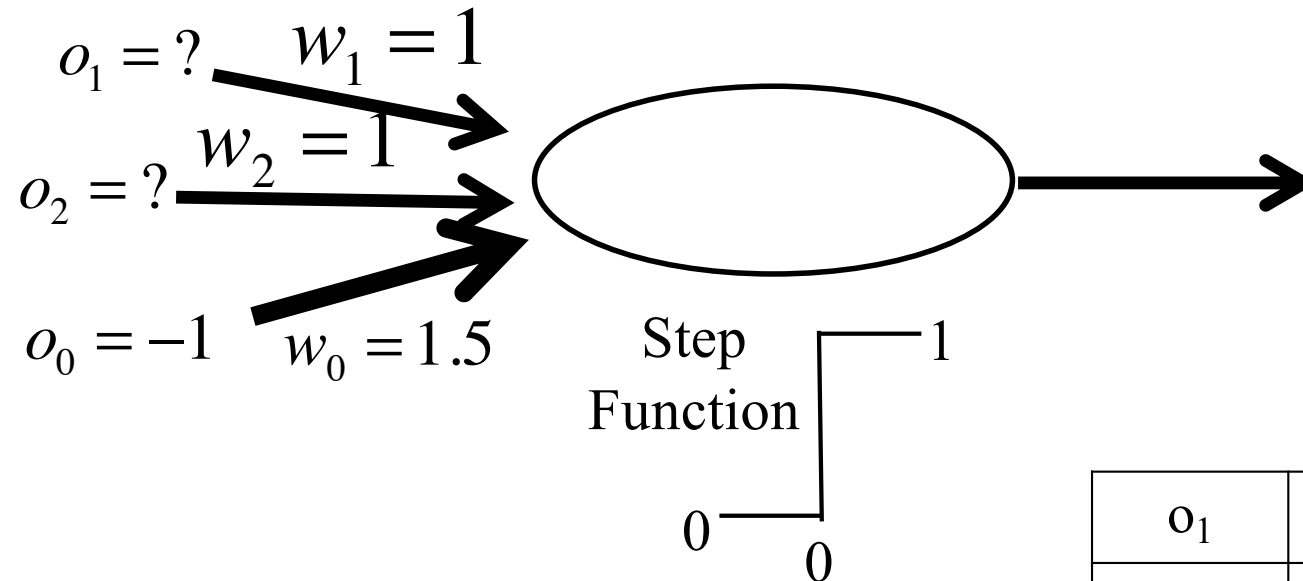
# McCulloch-Pitts “Neuron”



$$1 + 1 - 1.5 = 0.5$$

$o_1$	$o_2$	output
0	0	0
1	0	0
0	1	0
1	1	1

# McCulloch-Pitts “Neuron”



This neuron  
computes the  
function “and.”

There are “or” and “not” neurons too.

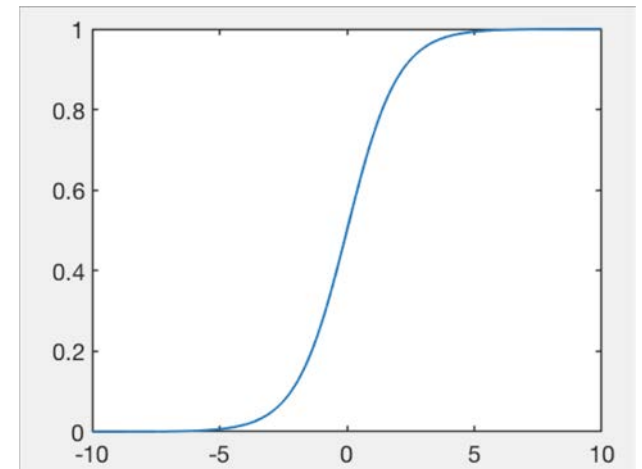
$o_1$	$o_2$	output
0	0	0
1	0	0
0	1	0
1	1	1

# McCulloch-Pitts “Neuron”

$$\phi\left(\sum_j w_{j,\text{me}} a_j\right) = 1 / (1 + e^{-x}) \quad \text{“Sigmoid”}$$

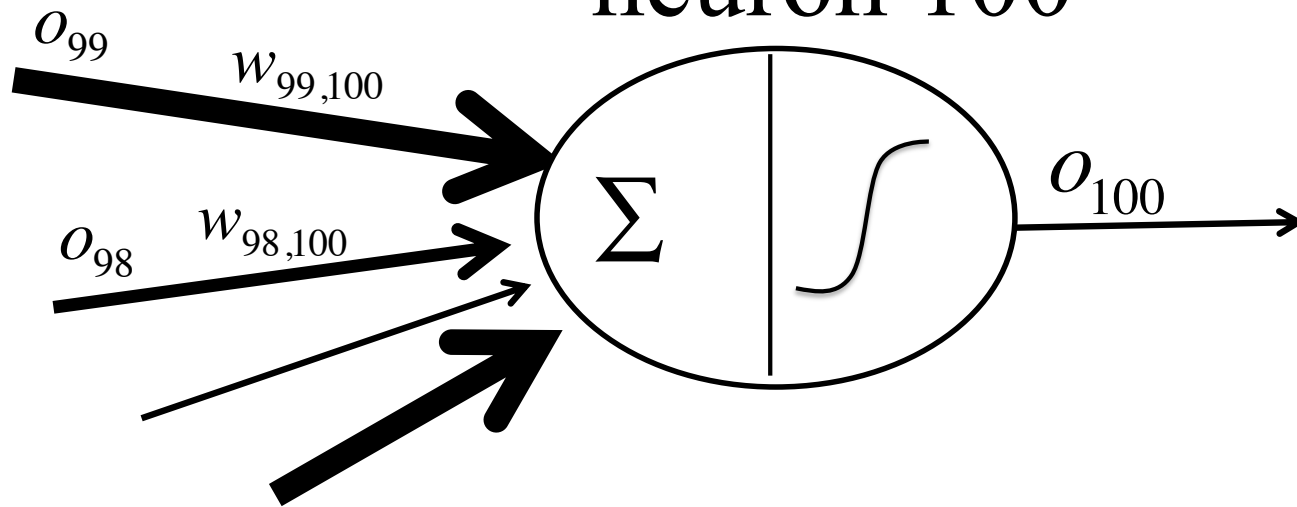


Activation  
Function



# Single Layer

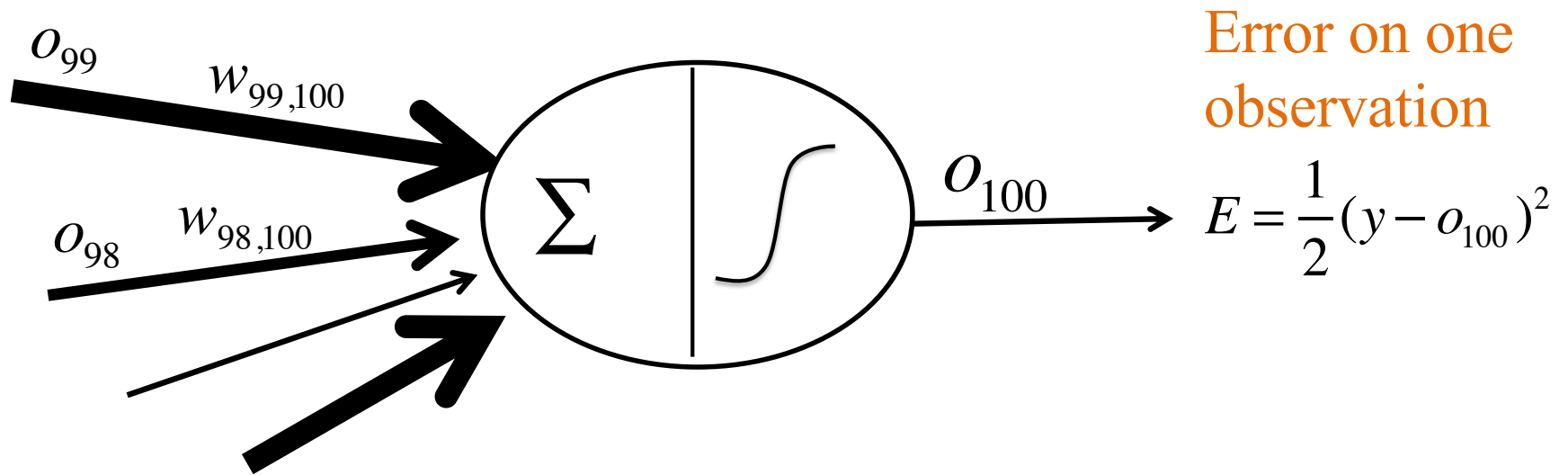
“neuron 100”



Error on one  
observation

$$E = \frac{1}{2}(y - o_{100})^2$$

# Single Layer



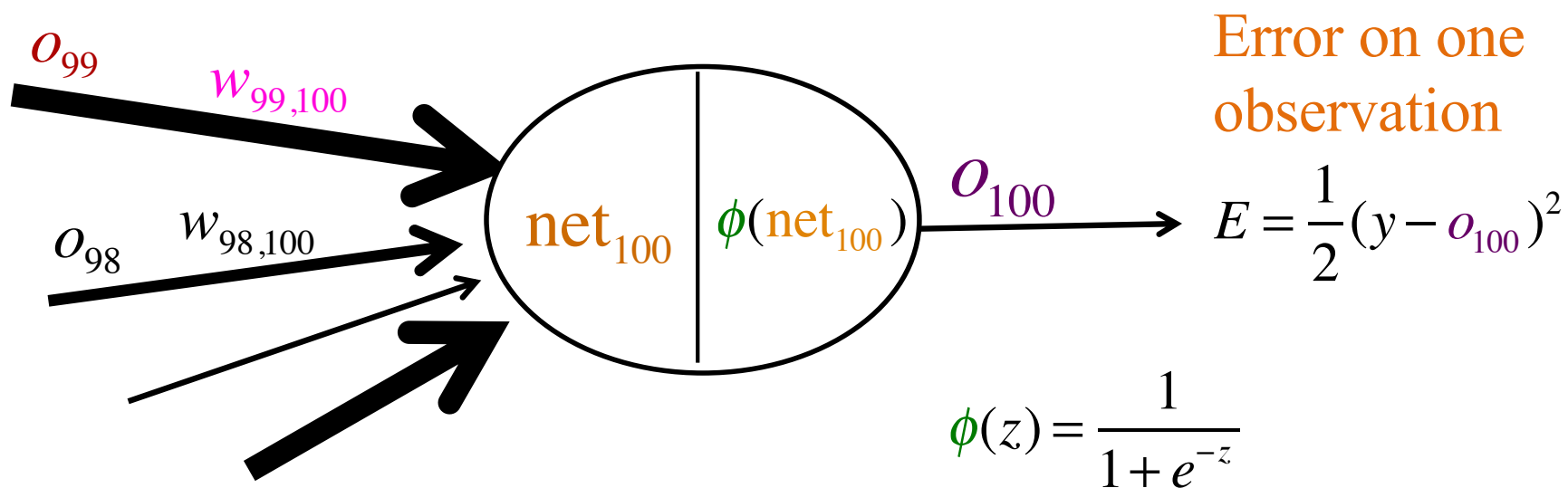
- In a brain, the synapses strengthen and weaken in order to learn.
- Say the same thing happens here.
- How should we set the weights in order to learn (reduce the error)?
- Minimize  $E$  with respect to the weights.

# Backpropagation

- An algorithm that trains the weights of a neural network
- Requires us to propagate information backwards through the network, then forwards, then backwards, then forwards, etc.
- Propagate backwards = chain rule from calculus.

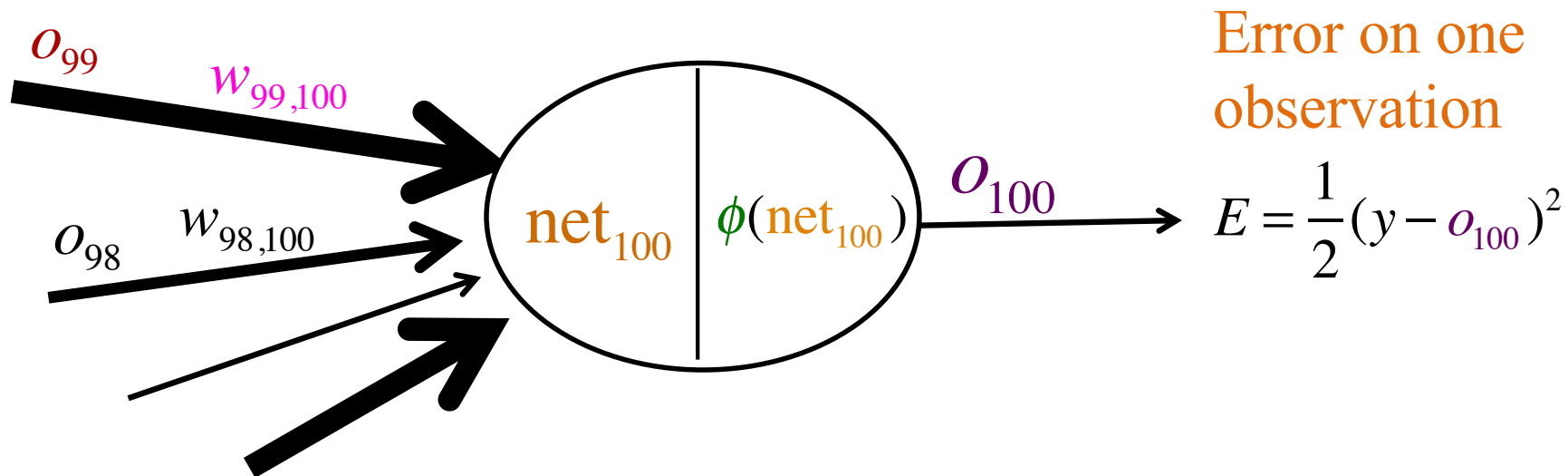


# Single Layer



$$\frac{dE}{dw_{99,100}} = \frac{dE}{dO_{100}} \frac{dO_{100}}{d\text{net}_{100}} \frac{d\text{net}_{100}}{dw_{99,100}}$$

# Single Layer



$$\frac{dE}{dw_{99,100}} = \frac{dE}{do_{100}} \frac{do_{100}}{d \text{net}_{100}} \frac{d \text{net}_{100}}{dw_{99,100}}$$

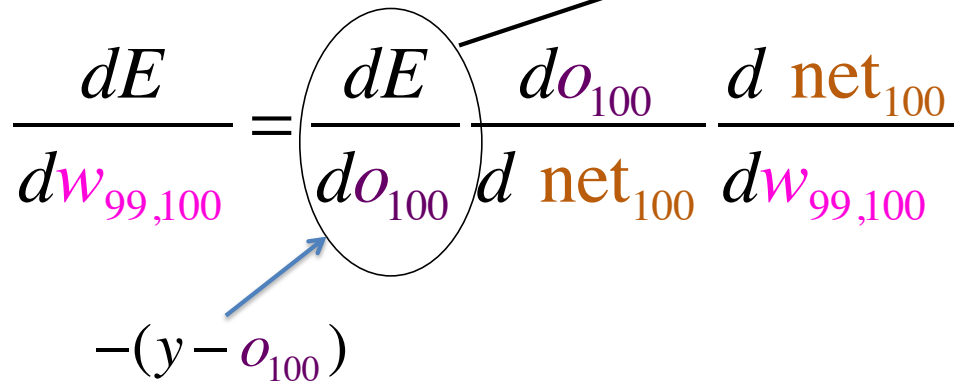
# Single Layer

Error on one  
observation

$$E = \frac{1}{2}(y - o_{100})^2$$

$$\frac{dE}{dw_{99,100}} = \frac{dE}{do_{100}} \frac{do_{100}}{d \text{net}_{100}} \frac{d \text{net}_{100}}{dw_{99,100}}$$

$-(y - o_{100})$



The picture can't be displayed.

# Single Layer

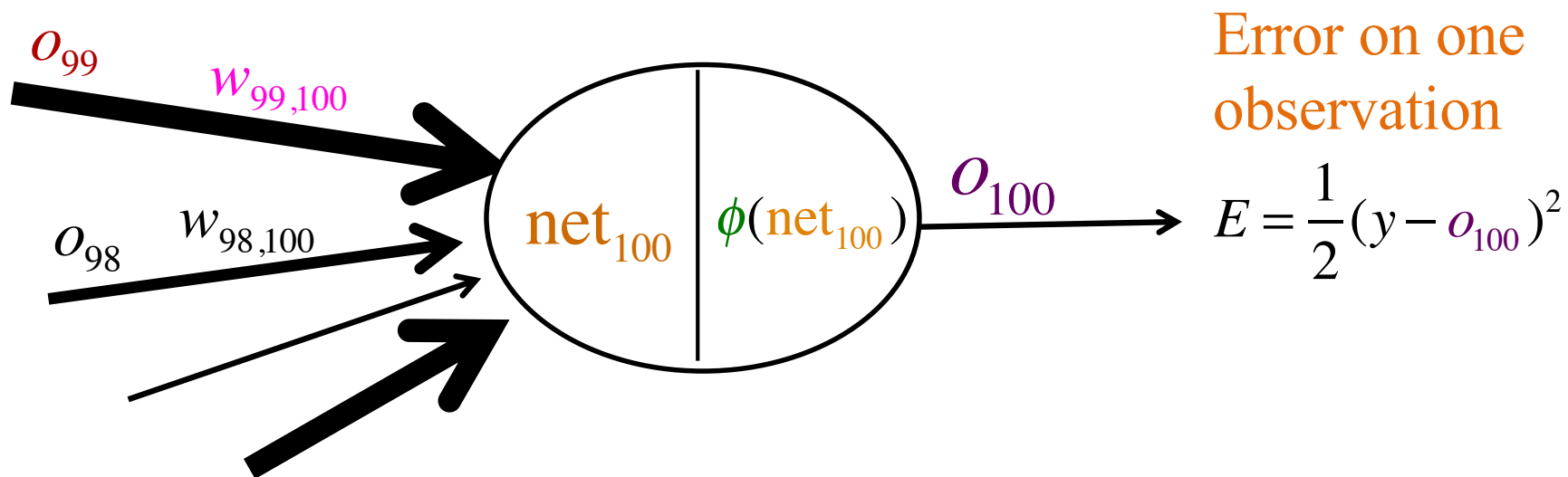
Error on one  
observation

$$E = \frac{1}{2}(y - o_{100})^2$$

$$\frac{dE}{d\mathbf{w}_{99,100}} = \left( \frac{dE}{do_{100}} \right) \frac{do_{100}}{d\mathbf{net}_{100}} \frac{d\mathbf{net}_{100}}{d\mathbf{w}_{99,100}}$$

$-(y - o_{100})$

# Single Layer



$$\frac{dE}{dw_{99,100}} = \frac{dE}{do_{100}} \left( \frac{do_{100}}{d \text{net}_{100}} \right) \frac{d \text{net}_{100}}{dw_{99,100}}$$

$-(y - o_{100})$

# Single Layer

$$\phi(\text{net}_{100}) \quad o_{100}$$

$$\frac{do_{100}}{d \text{net}_{100}} = \frac{d\phi(\text{net}_{100})}{d \text{net}_{100}} = \phi'(\text{net}_{100}) = \phi(\text{net}_{100})(1 - \phi(\text{net}_{100})) = o_{100}(1 - o_{100})$$

$$\frac{dE}{dw_{99,100}} = \frac{dE}{do_{100}} \left( \frac{do_{100}}{d \text{net}_{100}} \right) \frac{d \text{net}_{100}}{dw_{99,100}}$$

$$-(y - o_{100})$$

# Single Layer

$$\phi(\text{net}_{100}) \quad o_{100}$$

$$\frac{do_{100}}{d \text{net}_{100}} = \frac{d\phi(\text{net}_{100})}{d \text{net}_{100}} = \phi'(\text{net}_{100}) = \phi(\text{net}_{100})(1 - \phi(\text{net}_{100})) = o_{100}(1 - o_{100})$$


$$\frac{dE}{d w_{99,100}} = \frac{dE}{do_{100}} \left( \frac{do_{100}}{d \text{net}_{100}} \right) \frac{d \text{net}_{100}}{d w_{99,100}}$$

$-(y - o_{100})$        $o_{100}(1 - o_{100})$

# Single Layer

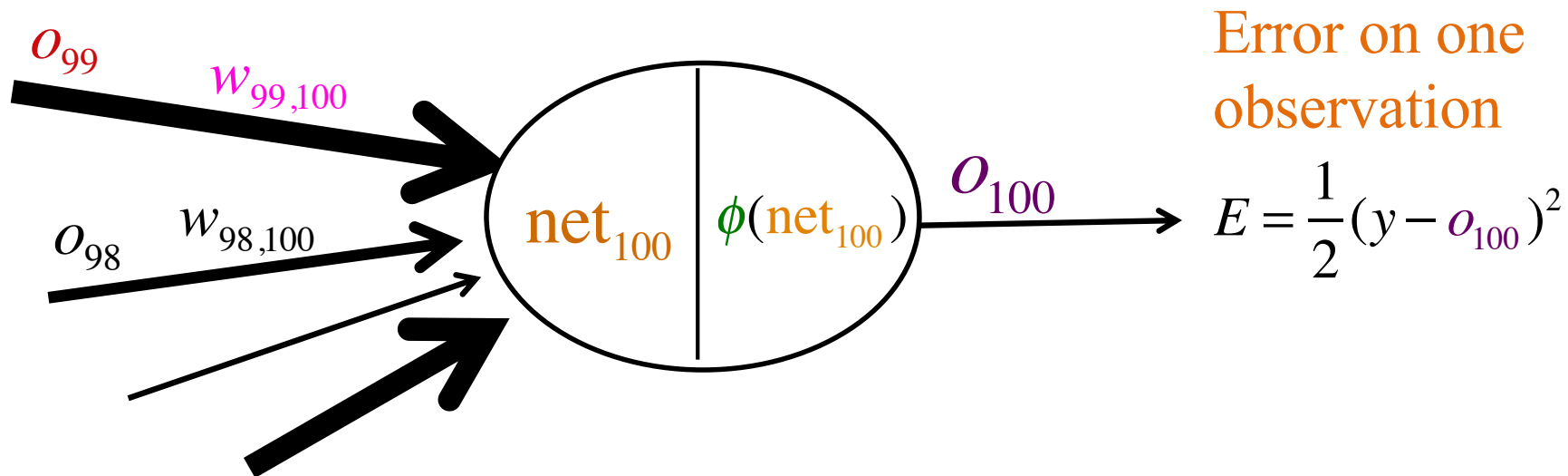
$$\frac{dE}{d\mathbf{w}_{99,100}} = \frac{dE}{d\mathbf{o}_{100}} \frac{d\mathbf{o}_{100}}{d\mathbf{net}_{100}} \frac{d\mathbf{net}_{100}}{d\mathbf{w}_{99,100}}$$

$-(y - \mathbf{o}_{100})$        $\mathbf{o}_{100}(1 - \mathbf{o}_{100})$





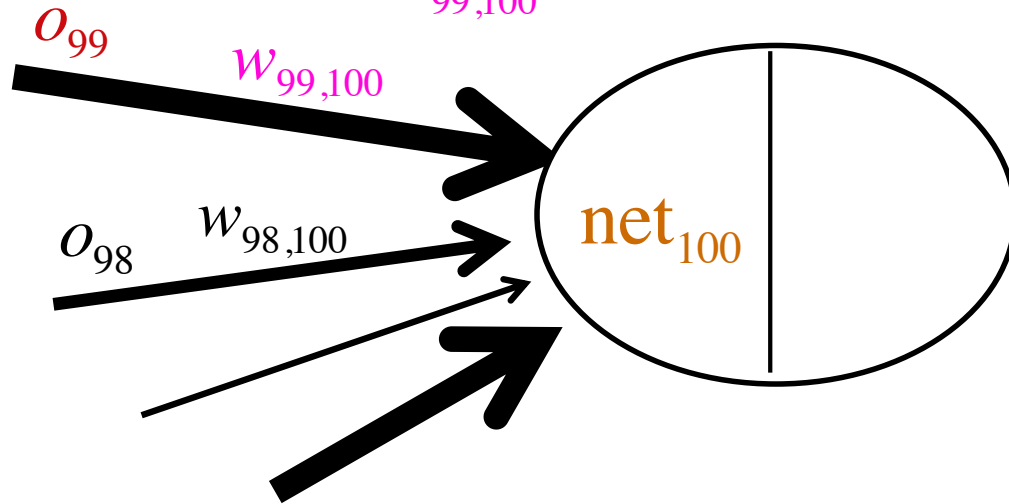
# Single Layer



$$\frac{dE}{dw_{99,100}} = \frac{dE}{dO_{100}} \frac{dO_{100}}{d\text{net}_{100}} \frac{d\text{net}_{100}}{dw_{99,100}}$$

$\frac{dE}{dO_{100}} = -(y - O_{100})$  (indicated by a blue arrow pointing to  $dO_{100}$ )  
 $\frac{d\text{net}_{100}}{dw_{99,100}} = O_{100}(1 - O_{100})$  (indicated by a blue arrow pointing to  $d\text{net}_{100}$ )


$$\frac{d \text{net}_{100}}{dw_{99,100}} = \frac{d (w_{99,100} o_{99} + w_{98,100} o_{98} + w_{97,100} o_{97} + \dots)}{dw_{99,100}} = o_{99}$$




$$\frac{dE}{dw_{99,100}} = \frac{dE}{do_{100}} \frac{do_{100}}{d \text{net}_{100}} \frac{d \text{net}_{100}}{dw_{99,100}}$$

Diagram illustrating the backpropagation of error gradients through the node  $\text{net}_{100}$ . The node is represented by a circle. Three blue arrows point towards the node from below, representing the error gradients for the weights  $w_{99,100}$ ,  $w_{98,100}$ , and  $w_{97,100}$ . The arrows are labeled with their respective error gradients:  $-(y - o_{100})$  for  $w_{99,100}$ ,  $o_{100}(1 - o_{100})$  for  $w_{98,100}$ , and  $o_{99}$  for  $w_{97,100}$ .


$$\frac{dE}{d\textcolor{violet}{w}_{99,100}} = \frac{dE}{d\textcolor{violet}{o}_{100}} \frac{d\textcolor{violet}{o}_{100}}{d\textcolor{brown}{net}_{100}} \frac{d\textcolor{brown}{net}_{100}}{d\textcolor{violet}{w}_{99,100}}$$



$-(y - \textcolor{violet}{o}_{100})$



$\textcolor{violet}{o}_{100} (1 - \textcolor{violet}{o}_{100})$



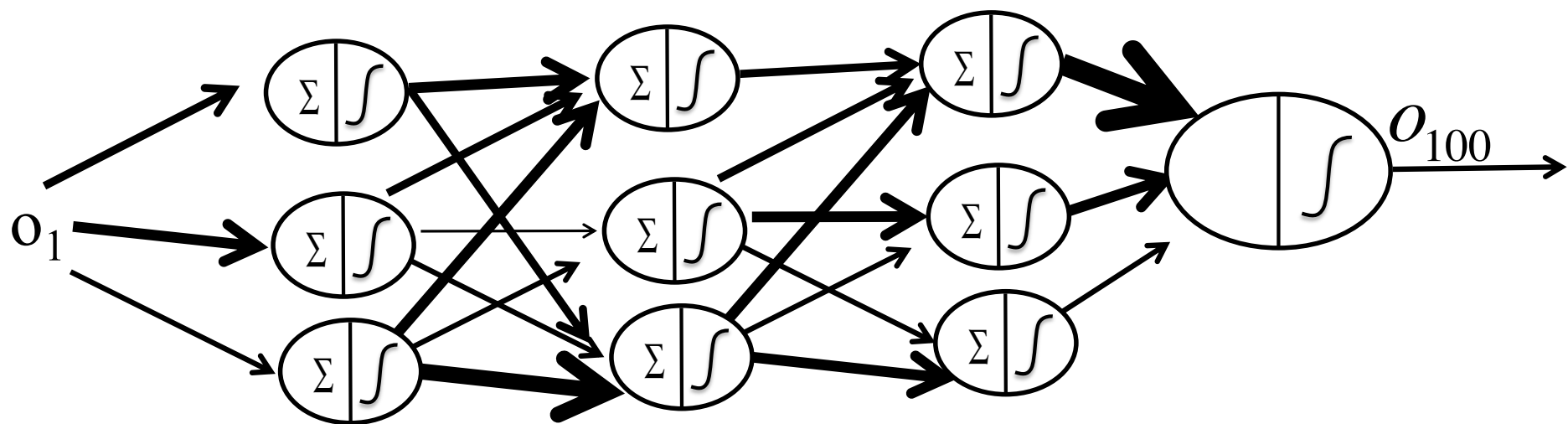
$\textcolor{red}{o}_{99}$

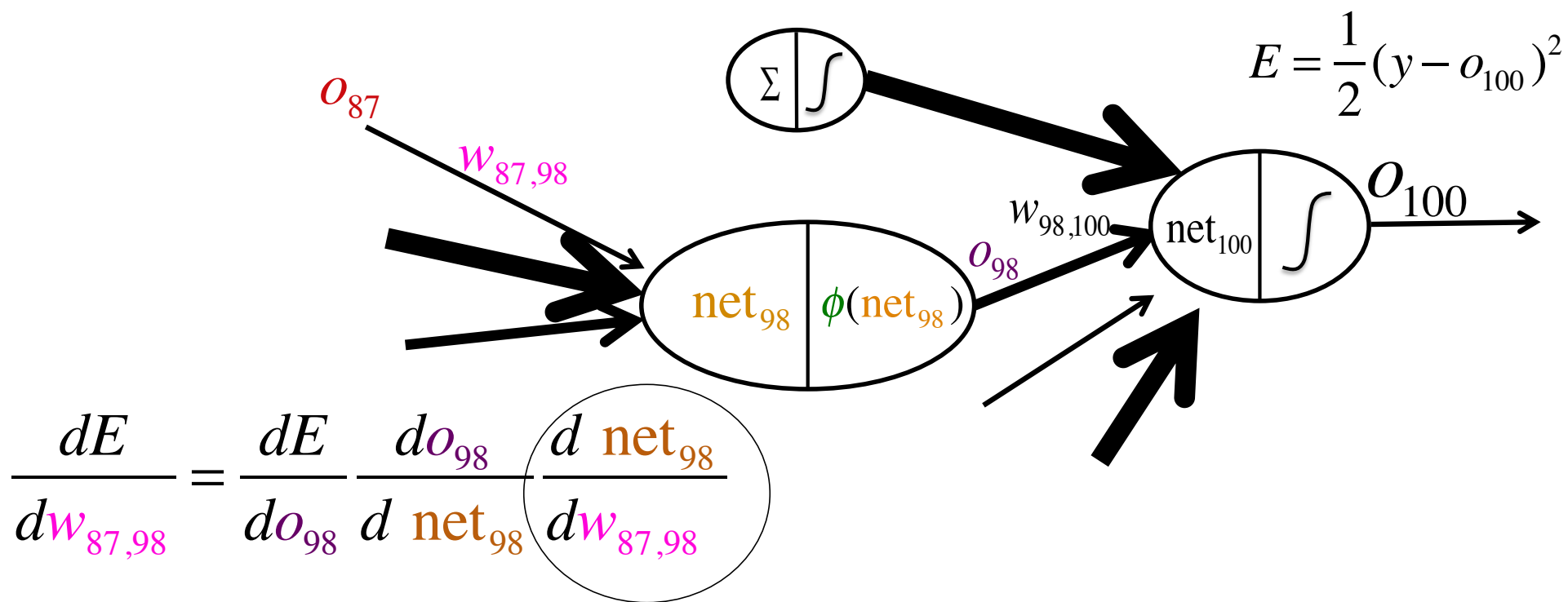
We will need this later – it depends only on node 100

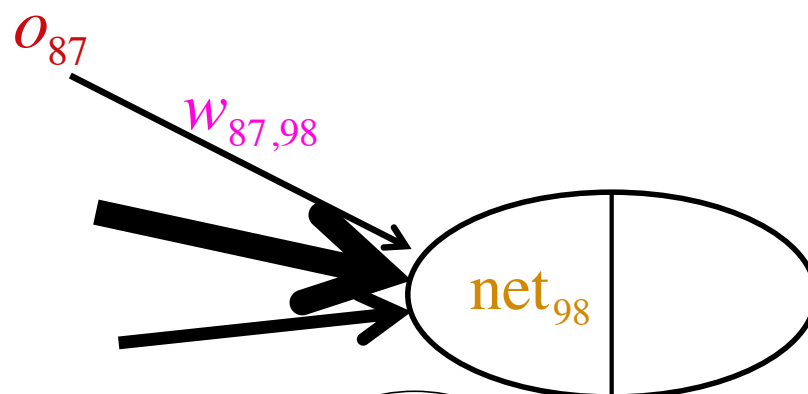
$$\begin{aligned}
 \frac{dE}{d\mathbf{w}_{99,100}} &= \frac{\delta_{100}}{\frac{dE}{do_{100}} \frac{d \mathbf{net}_{100}}{d \mathbf{w}_{99,100}}} \\
 &= (-(y - o_{100})) o_{100} (1 - o_{100}) o_{99}
 \end{aligned}$$

# Backpropagation

- Go one layer deeper.



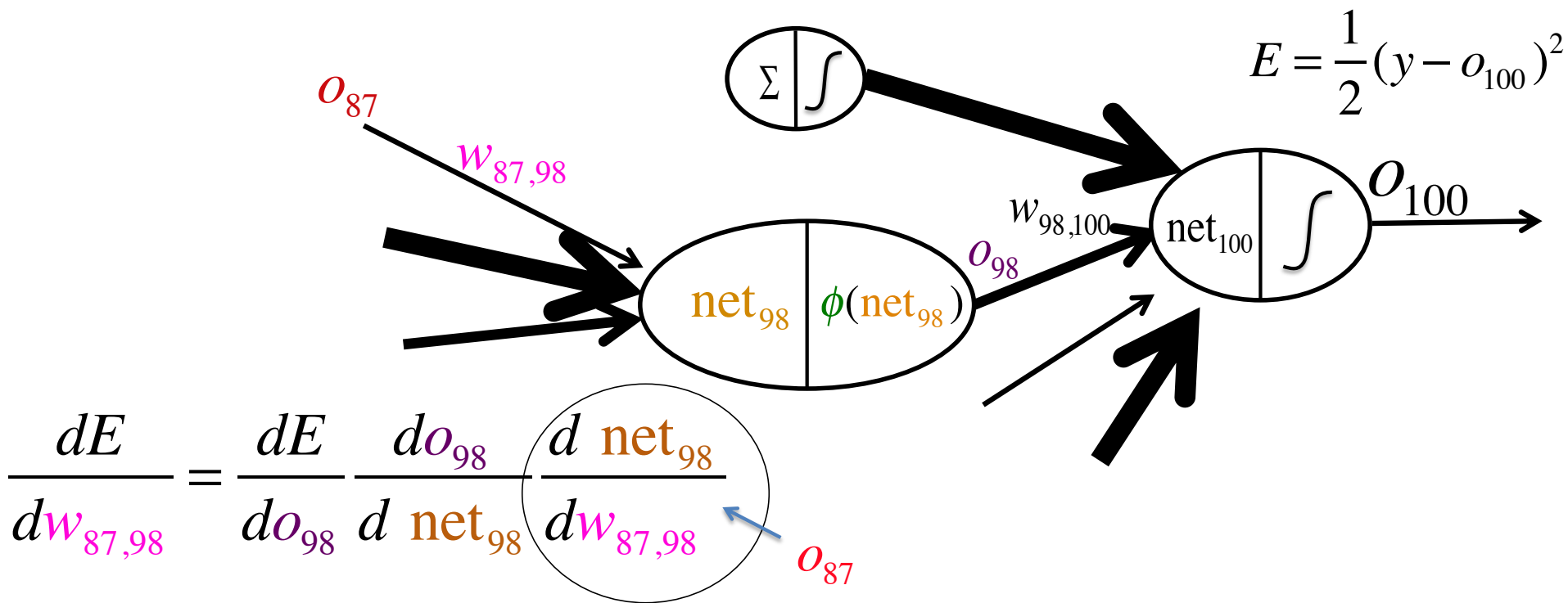


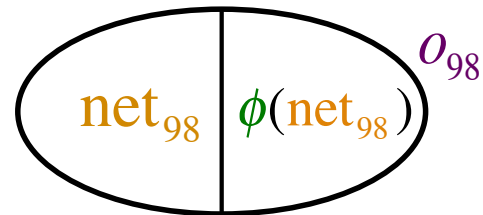


$$\frac{dE}{dw_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{d \text{net}_{98}} \frac{d \text{net}_{98}}{dw_{87,98}}$$

$$\frac{d \text{net}_{98}}{dw_{87,98}} = \frac{d (w_{87,98} O_{87} + w_{86,98} O_{86} + w_{85,98} O_{85} + \dots)}{dw_{87,98}} = O_{87}$$



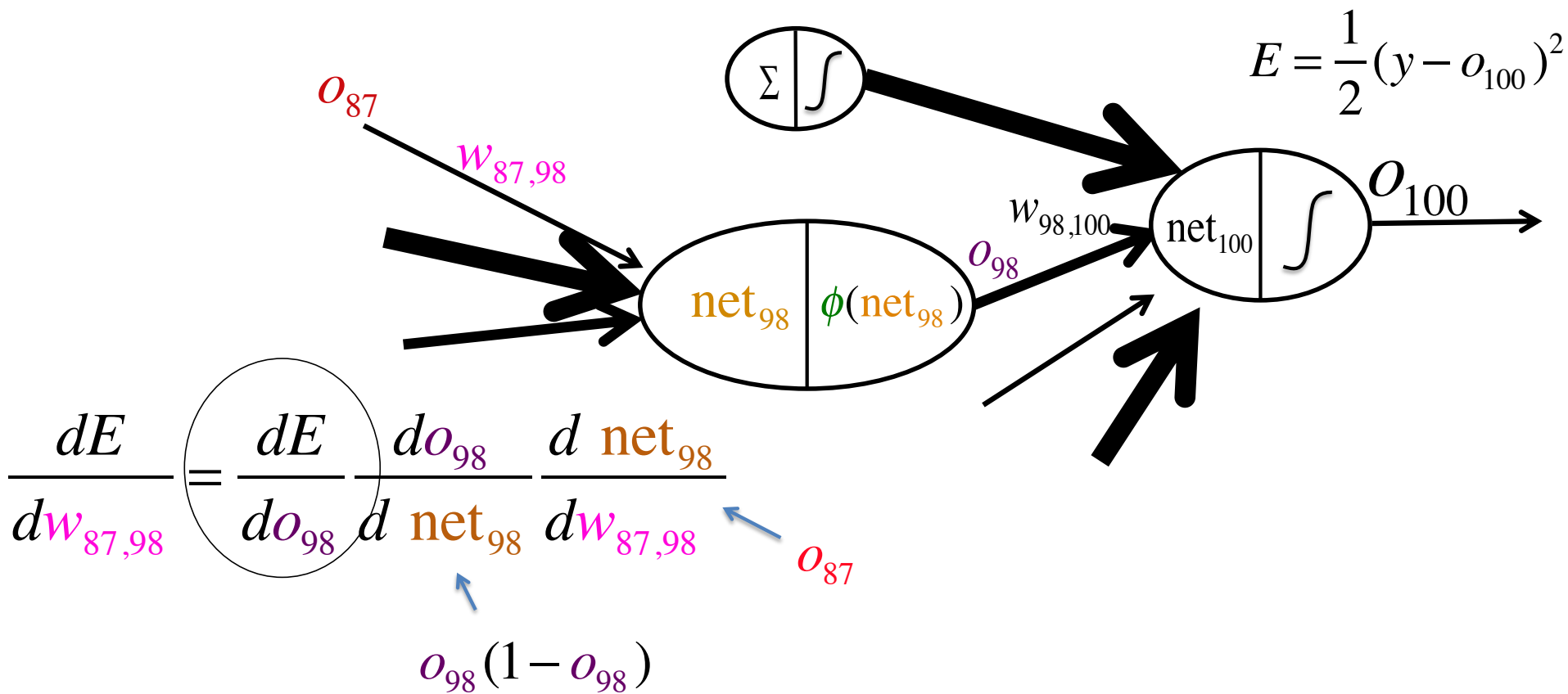


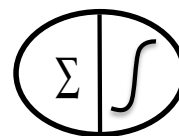


$$\frac{dE}{d\mathbf{w}_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{d\mathbf{net}_{98}} \frac{d\mathbf{net}_{98}}{d\mathbf{w}_{87,98}}$$

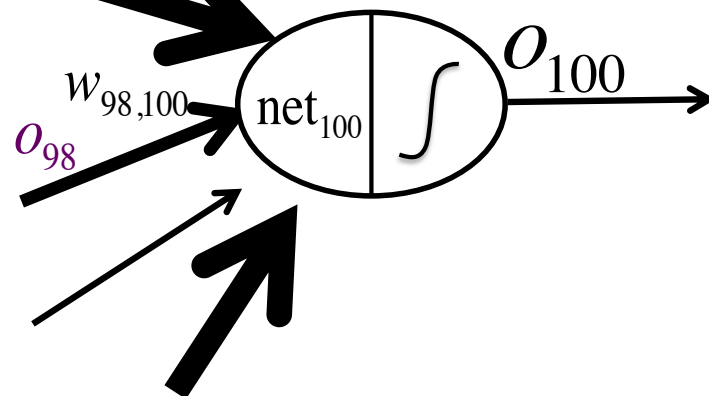
$\mathbf{o}_{87}$  ←

$$\frac{do_{98}}{d\mathbf{net}_{98}} = \frac{d\phi(\mathbf{net}_{98})}{d\mathbf{net}_{98}} = \phi'(\mathbf{net}_{98}) = \phi(\mathbf{net}_{98})(1 - \phi(\mathbf{net}_{98})) = o_{98}(1 - o_{98})$$





$$E = \frac{1}{2}(y - o_{100})^2$$



$$\frac{dE}{dw_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{d \text{net}_{98}} \frac{d \text{net}_{98}}{dw_{87,98}}$$

$o_{98}(1 - o_{98})$

$o_{87}$

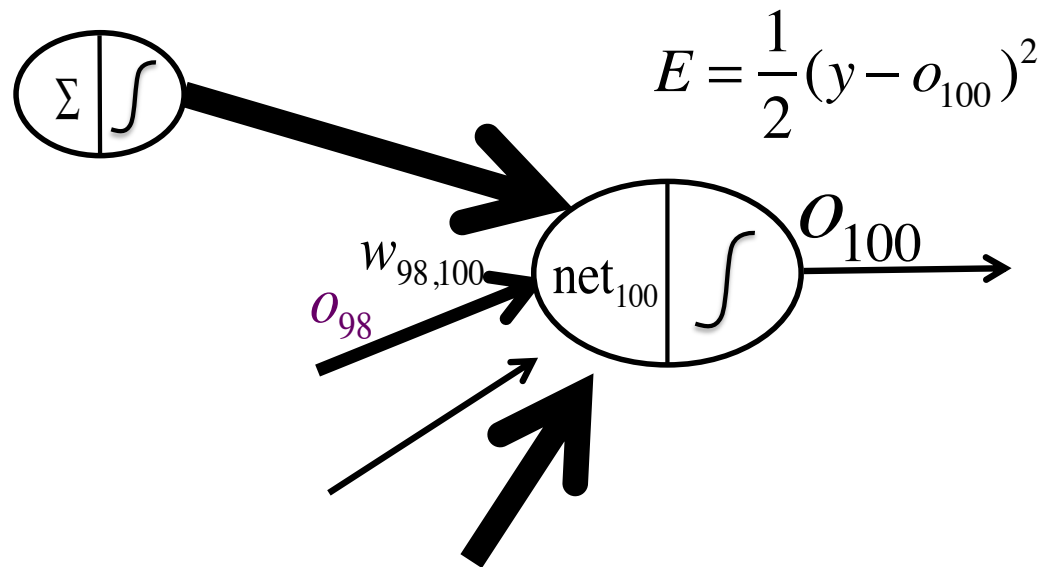
$$\frac{dE}{do_{98}} = \frac{dE}{d\text{net}_{100}} \frac{d\text{net}_{100}}{do_{98}}$$

$\delta_{100}$

$$\frac{dE}{dw_{87,98}} = \frac{dE}{do_{98}} \frac{do_{98}}{d\text{net}_{98}} \frac{d\text{net}_{98}}{dw_{87,98}}$$

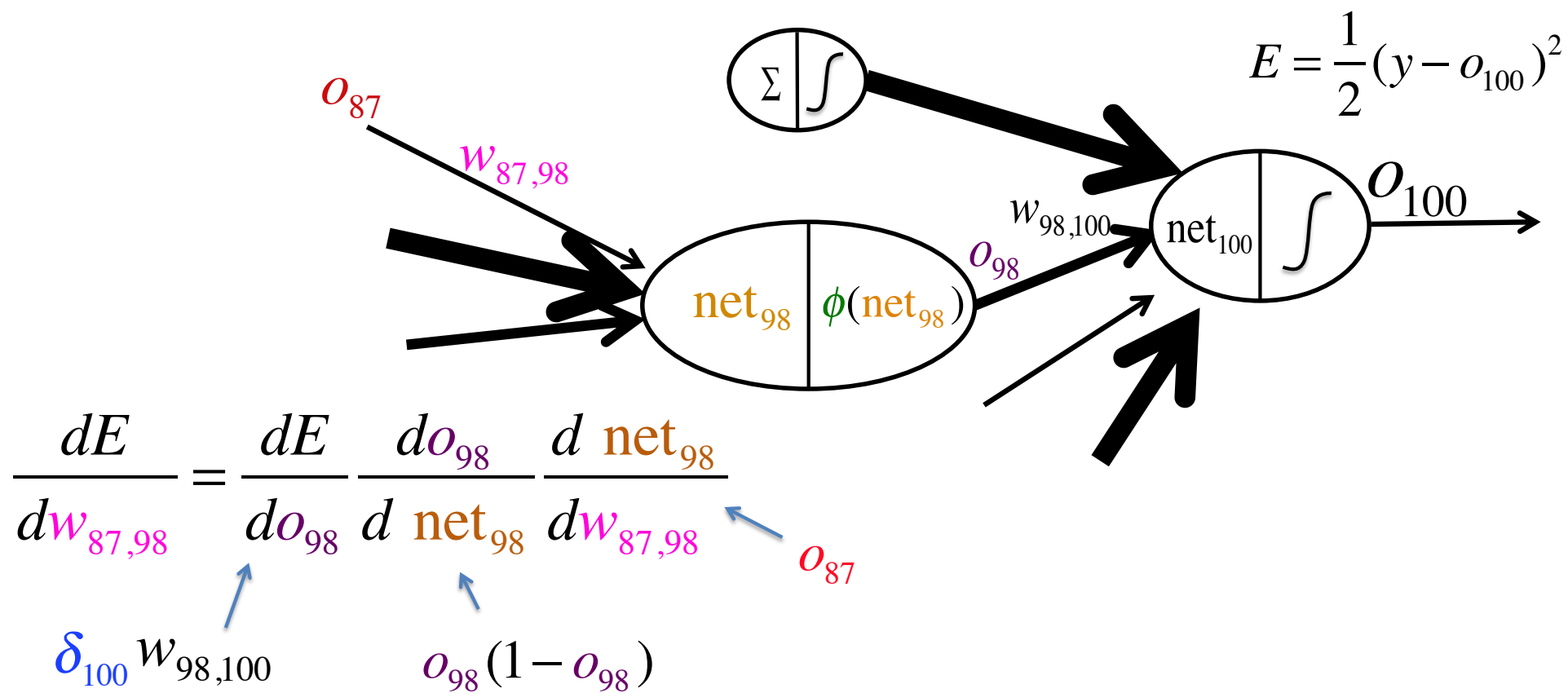
$o_{98}(1 - o_{98})$

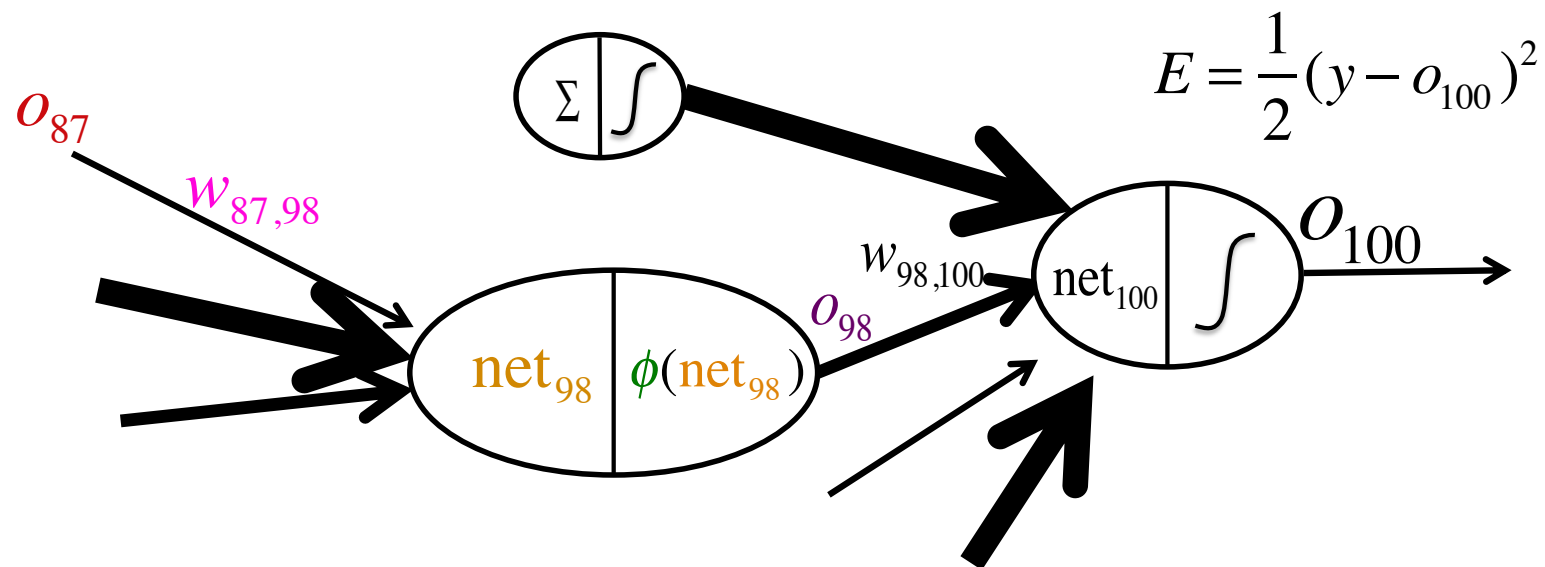
$o_{87}$



$$\begin{aligned}
 \frac{dE}{do_{98}} &= \frac{dE}{dnet_{100}} \frac{dnet_{100}}{do_{98}} = \frac{d(w_{99,100}o_{99} + w_{98,100}o_{98} + \dots)}{do_{98}} = w_{98,100} \\
 \frac{dE}{dw_{87,98}} &= \frac{dE}{do_{98}} \frac{do_{98}}{dnet_{98}} \frac{dnet_{98}}{dw_{87,98}} \\
 &\quad \delta_{100} \quad \quad \quad o_{98}(1 - o_{98}) \quad \quad \quad o_{87}
 \end{aligned}$$

The diagram illustrates the backpropagation of gradients through a sigmoid layer. The top equation shows the derivative of the error  $E$  with respect to the output  $o_{98}$ , which is the weight  $w_{98,100}$ . The bottom equation shows the derivative of  $E$  with respect to the weight  $w_{87,98}$ , which is the product of the derivative of  $E$  with respect to  $o_{98}$ , the derivative of  $o_{98}$  with respect to the net input  $net_{98}$ , and the derivative of  $net_{98}$  with respect to  $w_{87,98}$ . The diagram uses circles to group terms and arrows to show the flow of gradients.



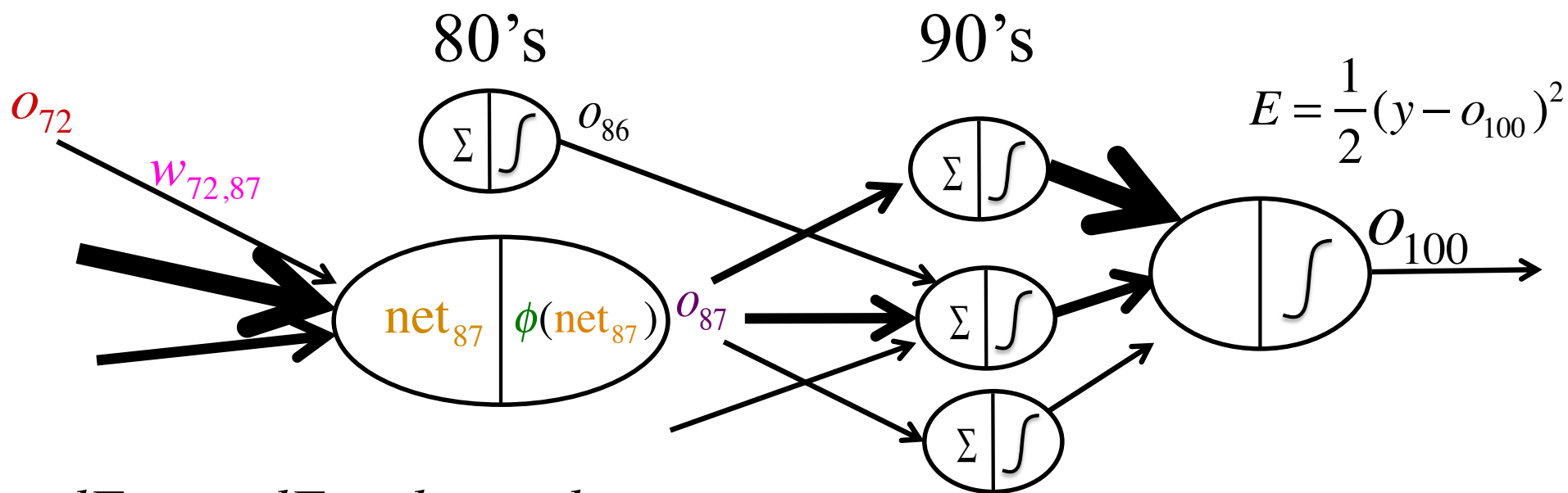


$$\frac{dE}{dw_{87,98}} = \delta_{100} w_{98,100} o_{98} (1 - o_{98}) o_{87}$$

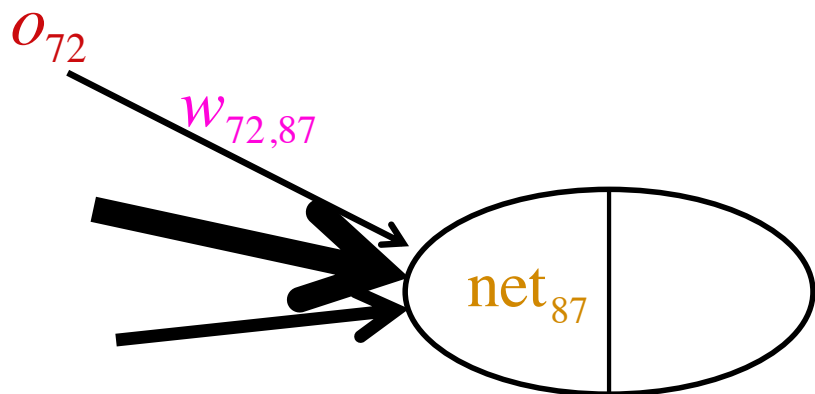


# Backpropagation

- Go even one layer deeper.
- Third time is a charm.

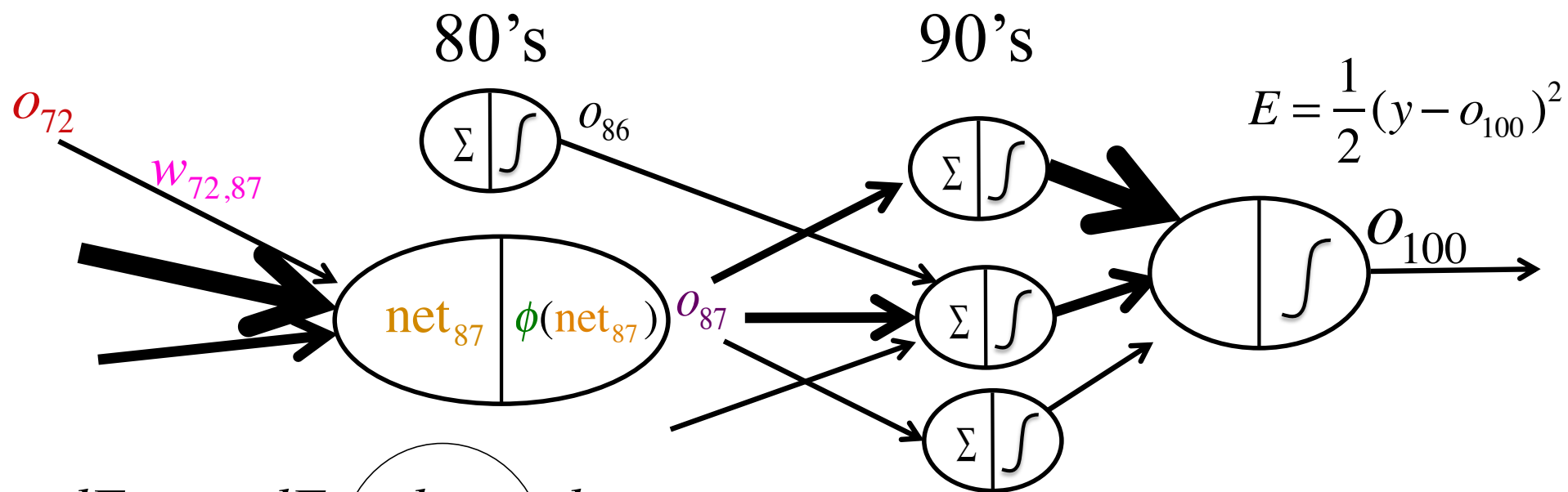


$$\frac{dE}{dw_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{d \text{net}_{87}} \frac{d \text{net}_{87}}{dw_{72,87}}$$



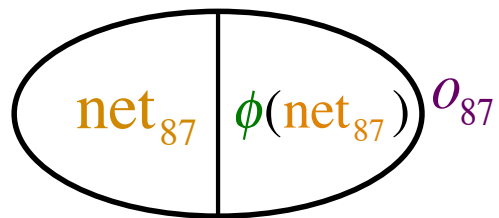
$$\frac{dE}{dw_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{d \text{net}_{87}} \frac{d \text{net}_{87}}{dw_{72,87}}$$

The term  $\frac{d \text{net}_{87}}{dw_{72,87}}$  in the equation is circled. A blue arrow points from a red label  $o_{72}$  to the denominator  $w_{72,87}$  of this term.



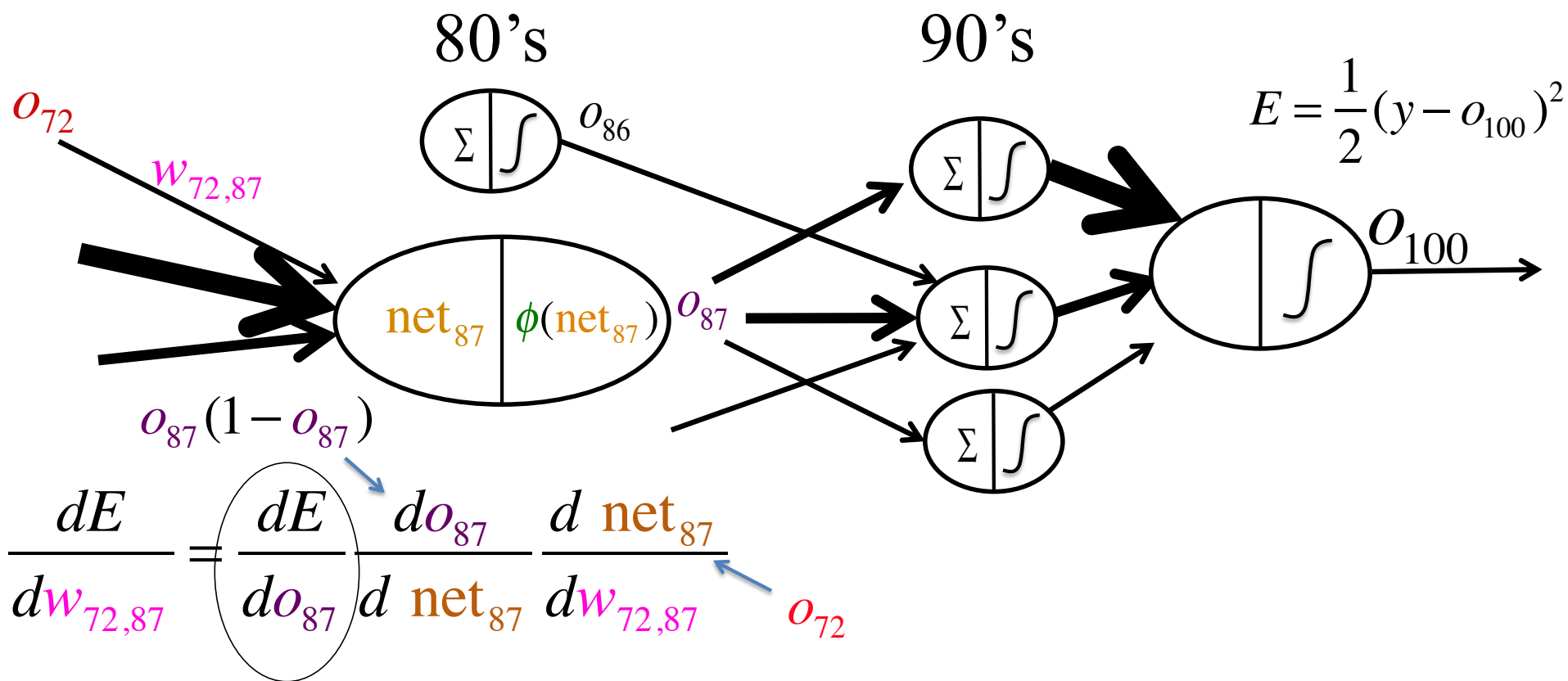
$$\frac{dE}{dw_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{d\text{net}_{87}} \frac{d\text{net}_{87}}{dw_{72,87}}$$

The term  $\frac{do_{87}}{d\text{net}_{87}}$  is circled in the original image. A blue arrow points from  $o_{72}$  (in red) to the term  $\frac{d\text{net}_{87}}{dw_{72,87}}$ .



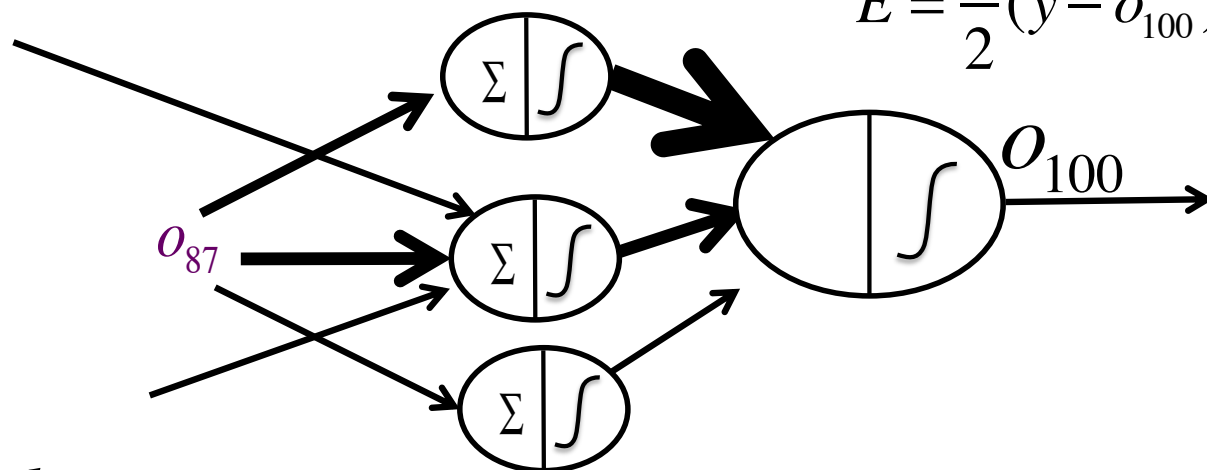
$$\frac{dE}{d\mathbf{w}_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{d\text{net}_{87}} \frac{d\text{net}_{87}}{d\mathbf{w}_{72,87}} \quad \text{with a blue arrow pointing to } \frac{d\text{net}_{87}}{d\mathbf{w}_{72,87}} \text{ and a red label } o_{72} \text{ next to it}$$

$$\frac{do_{87}}{d\text{net}_{87}} = \frac{d\phi(\text{net}_{87})}{d\text{net}_{87}} = \phi'(\text{net}_{87}) = \phi(\text{net}_{87})(1 - \phi(\text{net}_{87})) = o_{87}(1 - o_{87})$$



90's

$$E = \frac{1}{2}(y - o_{100})^2$$

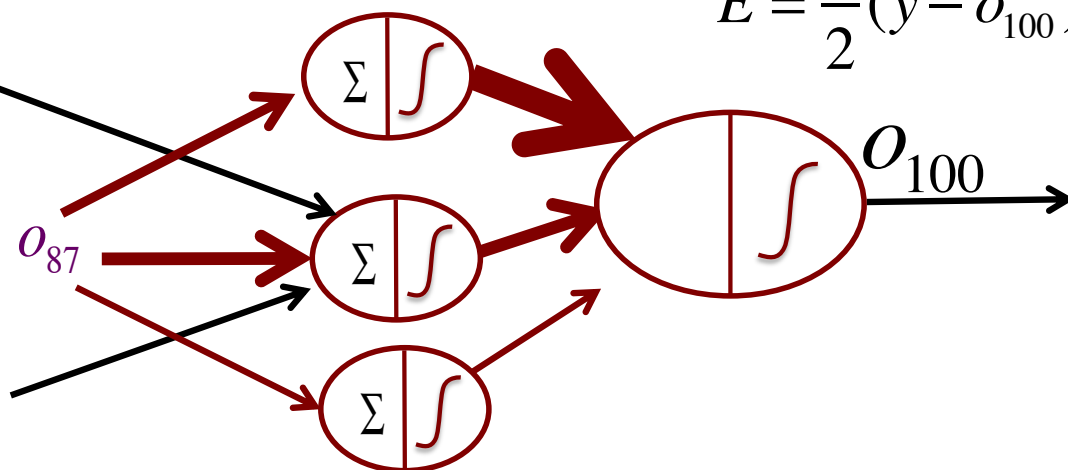


$$\frac{dE}{dw_{72,87}} = \left( \frac{dE}{do_{87}} \right) \frac{do_{87}}{d \text{net}_{87}} \frac{d \text{net}_{87}}{dw_{72,87}}$$

The equation shows the backpropagation of error gradients. The term  $\frac{dE}{do_{87}}$  is circled. A blue arrow points from the expression  $o_{87}(1 - o_{87})$  to  $\frac{dE}{do_{87}}$ . Another blue arrow points from the red label  $o_{72}$  to  $\frac{d \text{net}_{87}}{dw_{72,87}}$ .

90's

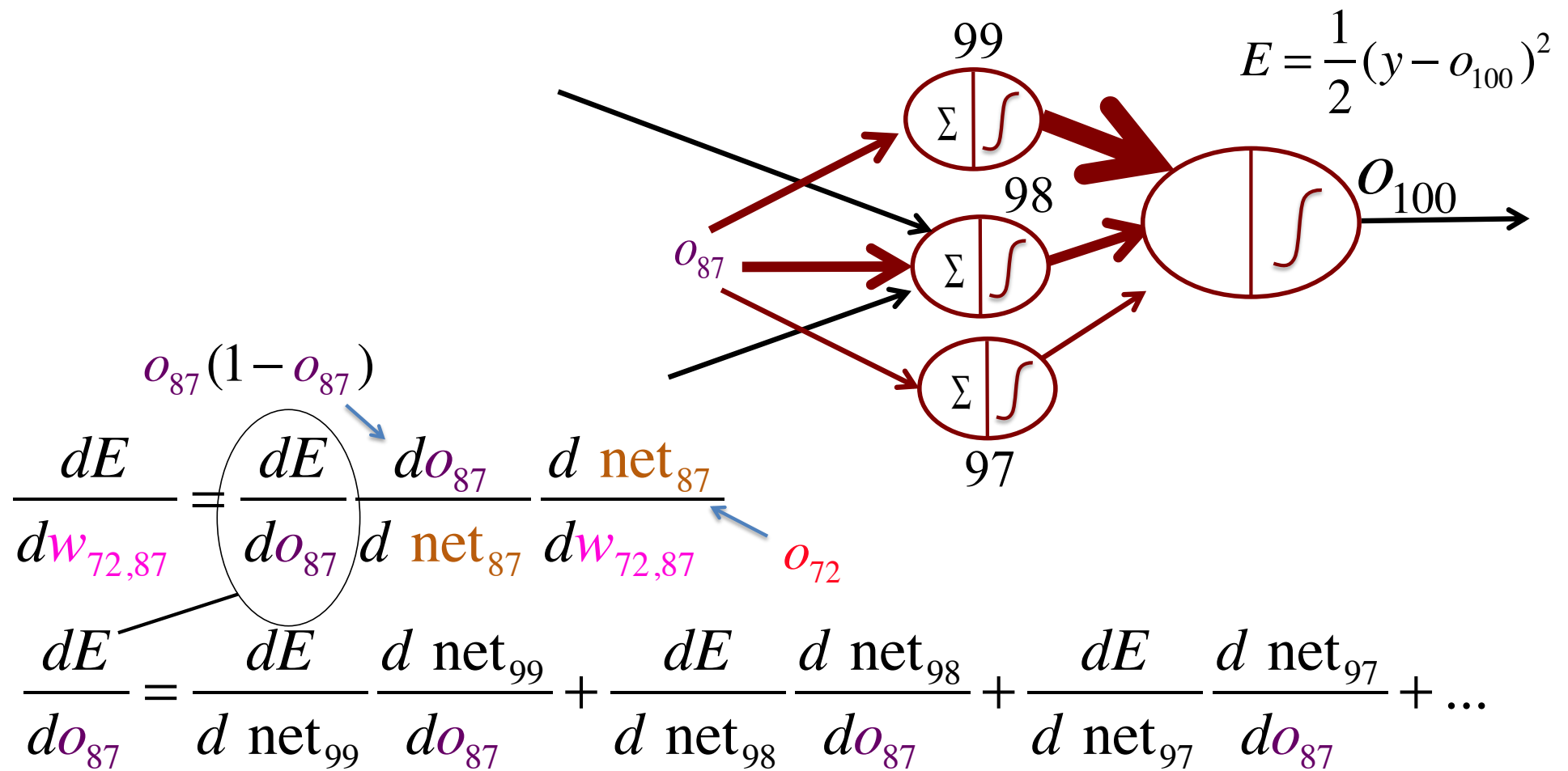
$$E = \frac{1}{2}(y - o_{100})^2$$



$$\frac{dE}{dw_{72,87}} = \left( \frac{dE}{do_{87}} \right) \frac{do_{87}}{d \text{net}_{87}} \frac{d \text{net}_{87}}{dw_{72,87}}$$

The equation shows the backpropagation of error gradients. The term  $\frac{dE}{do_{87}}$  is circled. A blue arrow points from  $o_{87}(1 - o_{87})$  to  $\frac{dE}{do_{87}}$ . Another blue arrow points from  $o_{72}$  to  $\frac{d \text{net}_{87}}{dw_{72,87}}$ .





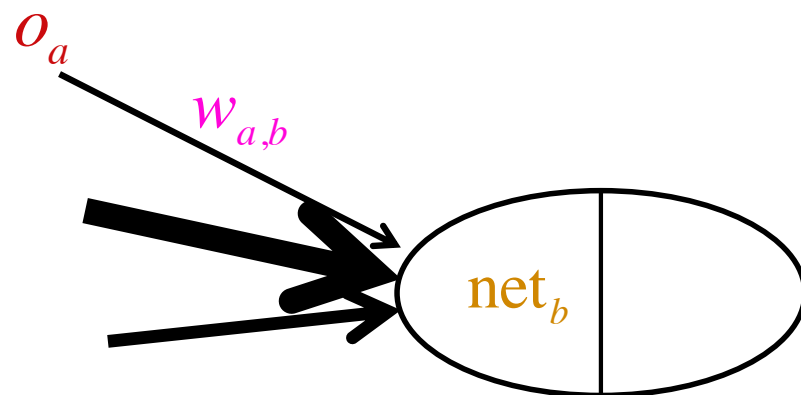
$$E = \frac{1}{2}(y - o_{100})^2$$

$$\frac{dE}{dw_{72,87}} = \frac{dE}{do_{87}} \frac{do_{87}}{d \text{net}_{87}} \frac{d \text{net}_{87}}{dw_{72,87}}$$

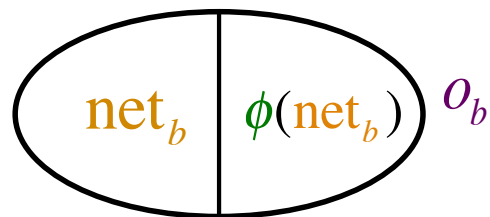
$$\frac{dE}{do_{87}} = \frac{dE}{d \text{net}_{99}} \frac{d \text{net}_{99}}{do_{87}} + \frac{dE}{d \text{net}_{98}} \frac{d \text{net}_{98}}{do_{87}} + \frac{dE}{d \text{net}_{97}} \frac{d \text{net}_{97}}{do_{87}} + \dots$$


$$= \delta_{99} w_{87,99} + \delta_{98} w_{87,98} + \delta_{97} w_{87,97} + \dots = \sum_{\ell \in L} \delta_{\ell} w_{87,\ell}$$

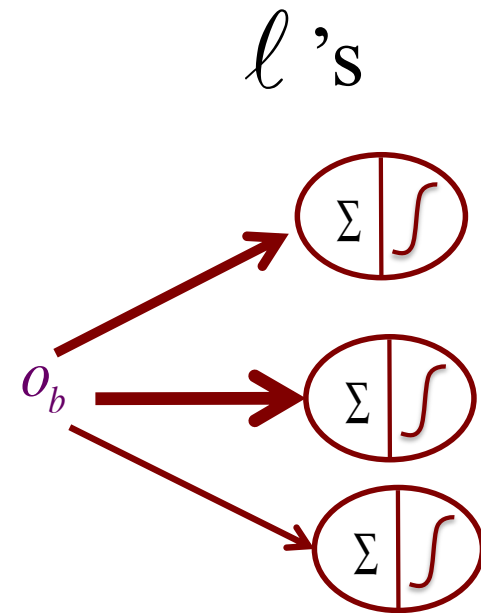
$$\begin{aligned}
 \frac{dE}{d\mathbf{w}_{a,b}} &= \frac{dE}{do_b} \frac{do_b}{d \mathbf{net}_b} \frac{d \mathbf{net}_b}{d\mathbf{w}_{a,b}} \\
 &= \frac{dE}{do_b} \frac{do_b}{d \mathbf{net}_b} o_a
 \end{aligned}$$



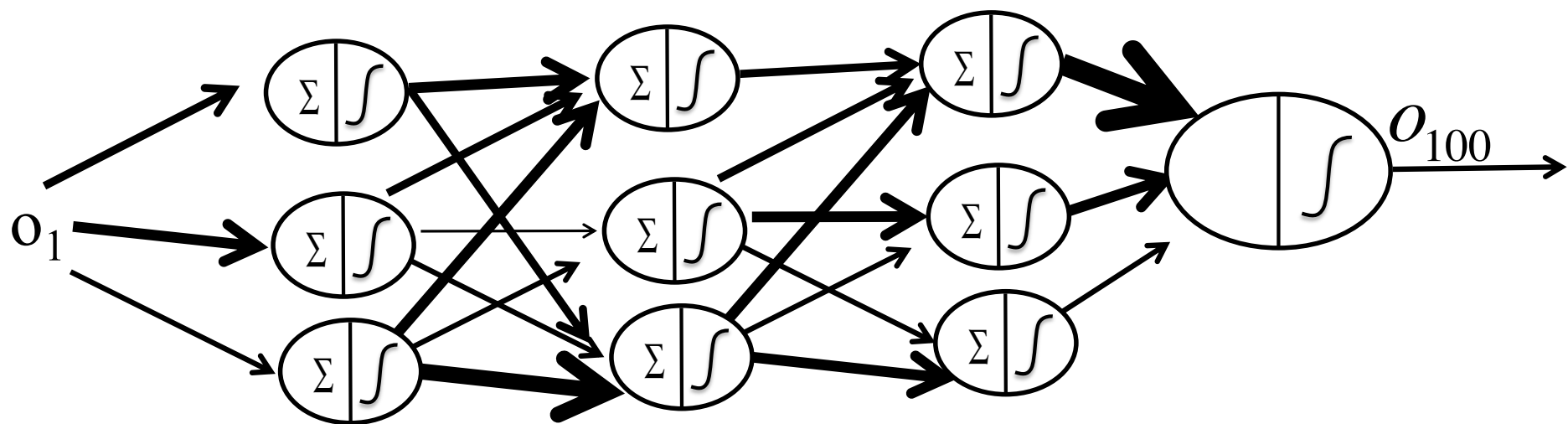
$$\begin{aligned}
 \frac{dE}{d\textcolor{violet}{w}_{a,b}} &= \frac{dE}{d\textcolor{violet}{o}_b} \frac{d\textcolor{violet}{o}_b}{d\textcolor{brown}{net}_b} \frac{d\textcolor{brown}{net}_b}{d\textcolor{violet}{w}_{a,b}} \\
 &= \frac{dE}{d\textcolor{violet}{o}_b} \textcolor{violet}{o}_b (1 - \textcolor{violet}{o}_b) \textcolor{violet}{o}_a
 \end{aligned}$$

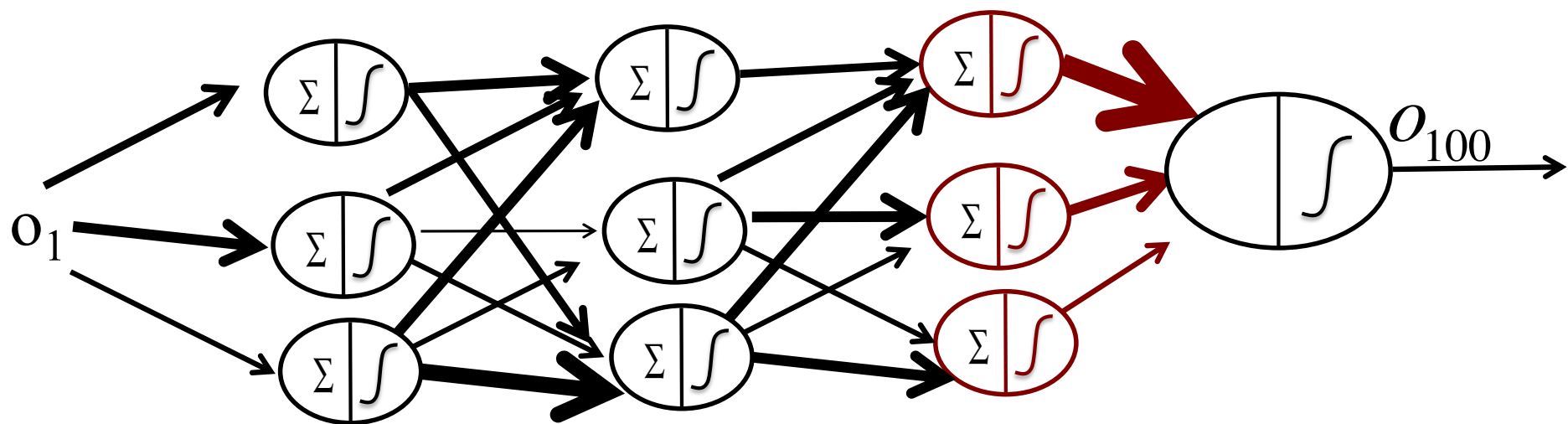


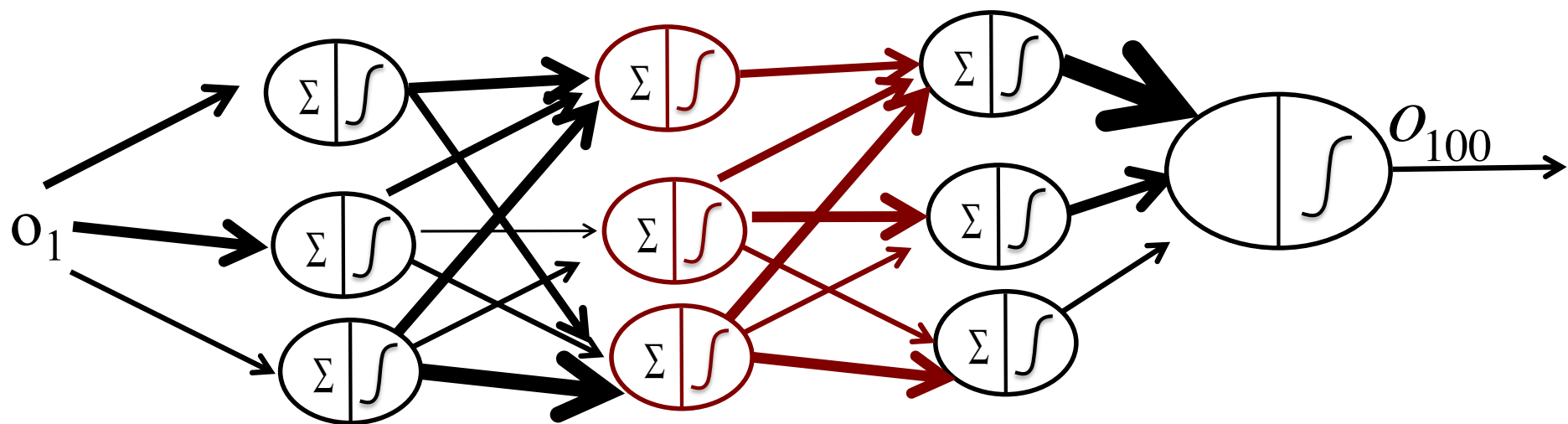
$$\begin{aligned}
 \frac{dE}{d\mathbf{w}_{a,b}} &= \frac{dE}{do_b} \frac{do_b}{d \text{net}_b} \frac{d \text{net}_b}{d\mathbf{w}_{a,b}} \\
 &= \left( \sum_{\ell \in L} \delta_\ell w_{b,\ell} \right) o_b (1 - o_b) o_a
 \end{aligned}$$




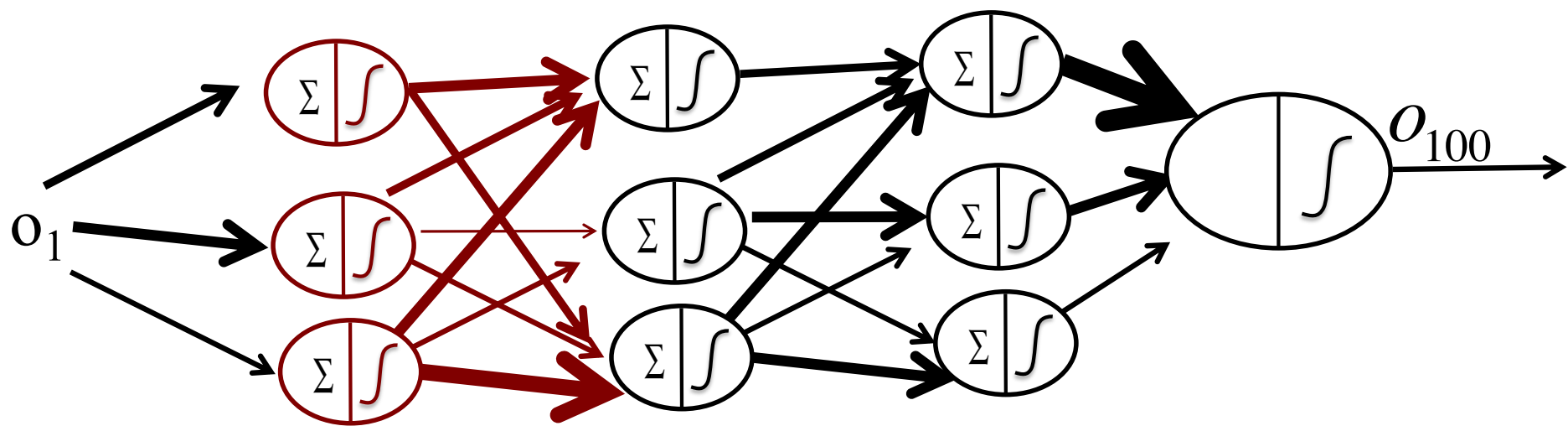
The  $\ell$ 's are downstream. We must have already computed all the  $\delta_\ell$ 's ahead of us to compute this.

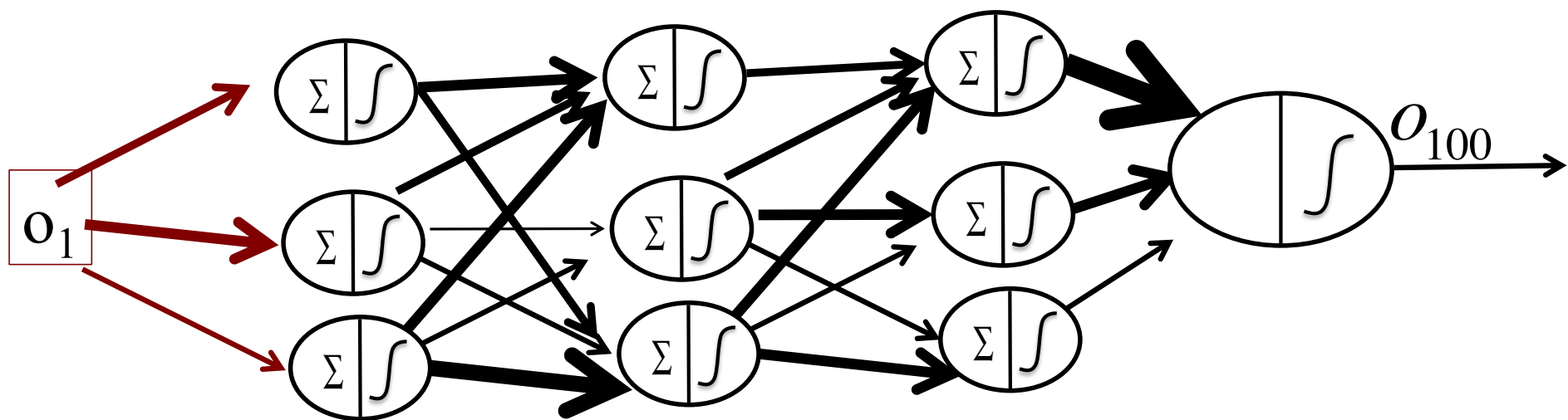


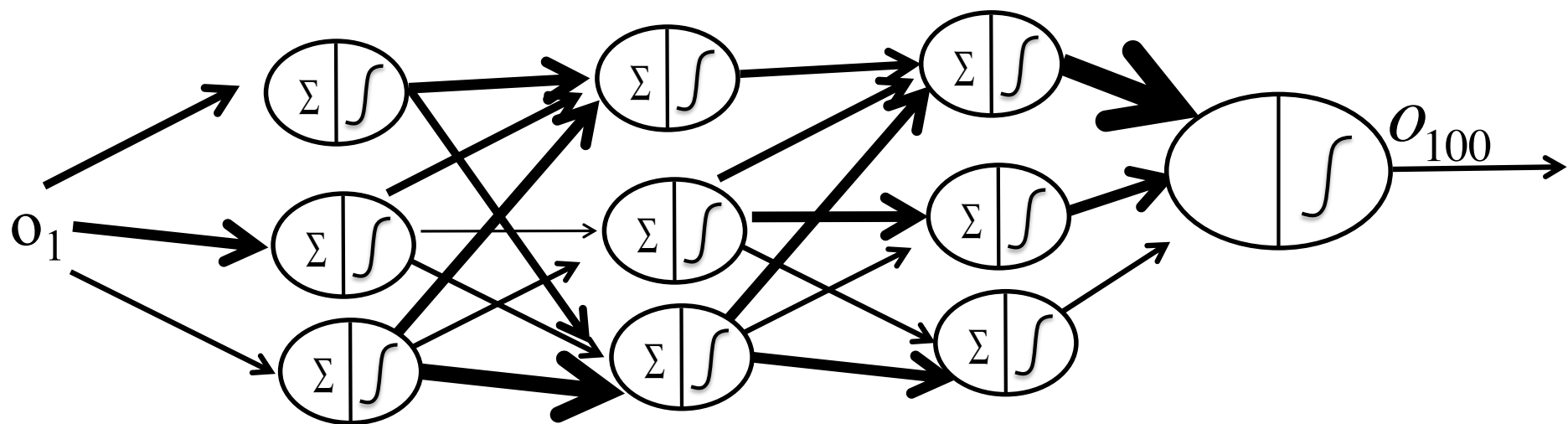










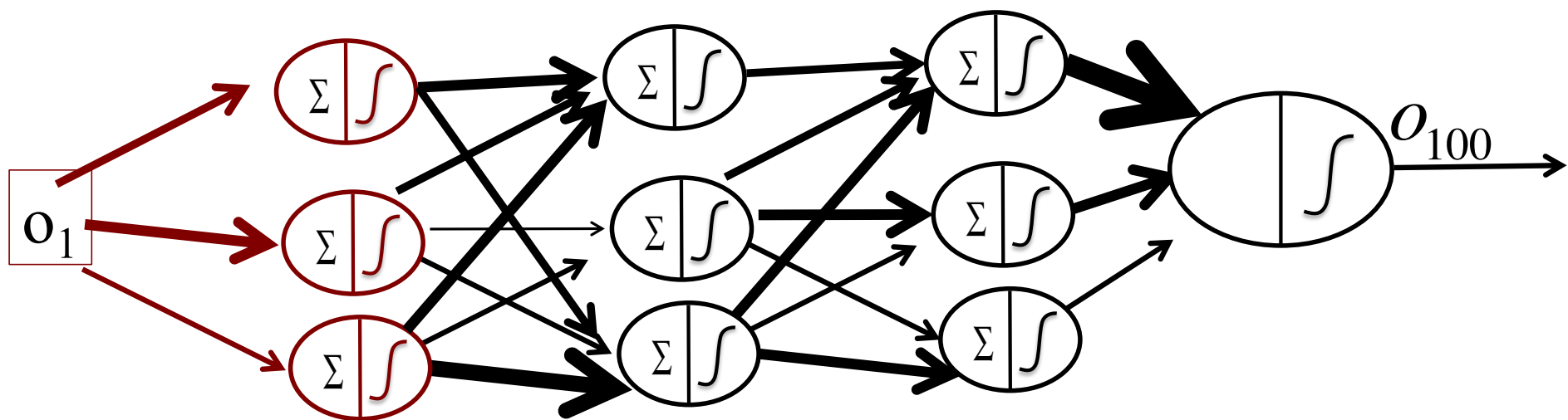


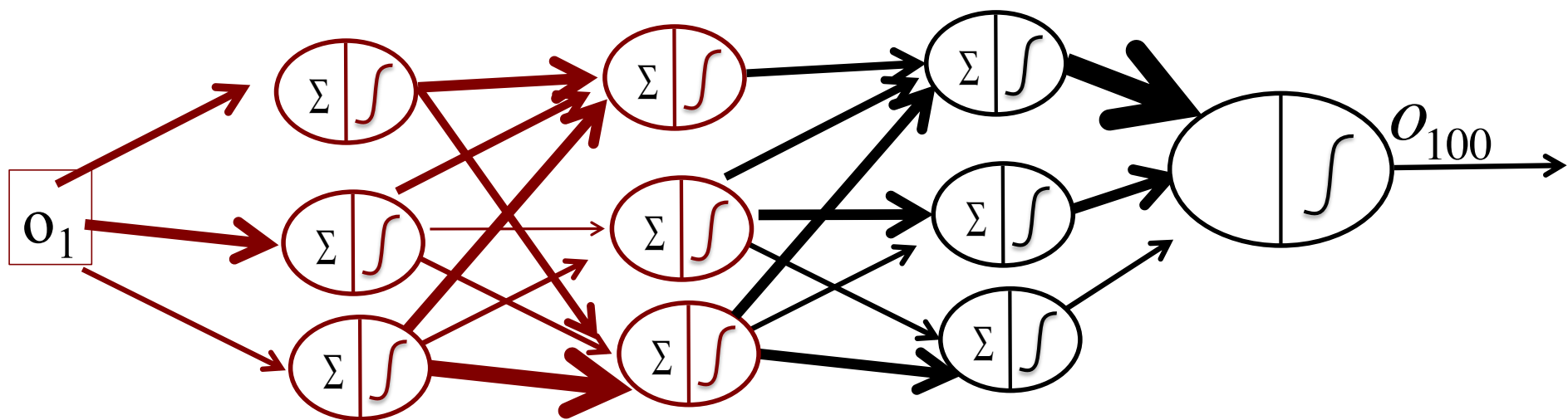
# Backpropagation

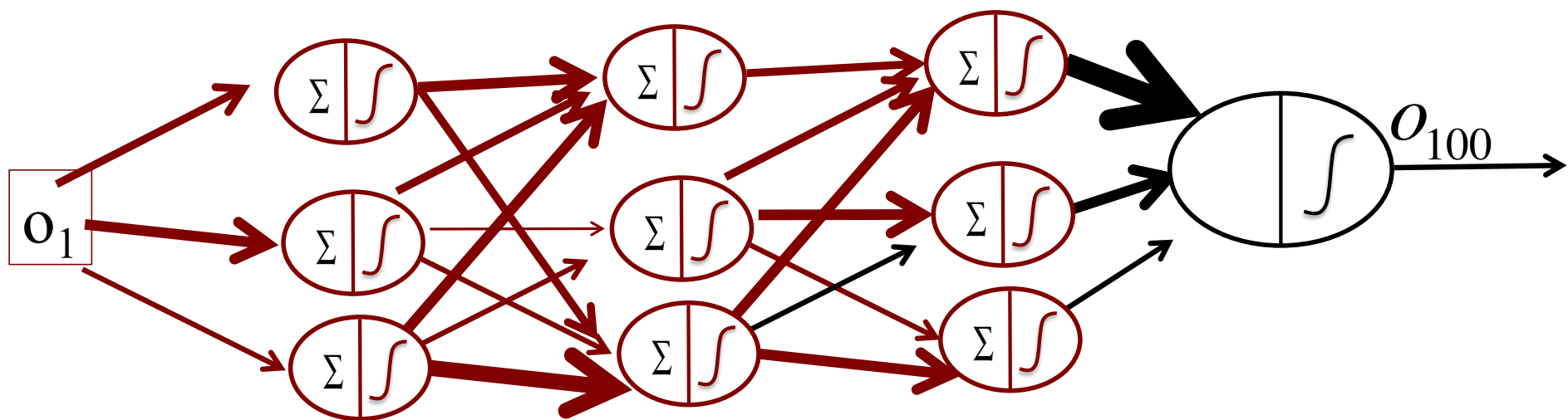
- Now we know how to compute  $\frac{dE}{dw_{a,b}}$  for all of the  $w_{a,b}$ 's.
- Let's do gradient descent.

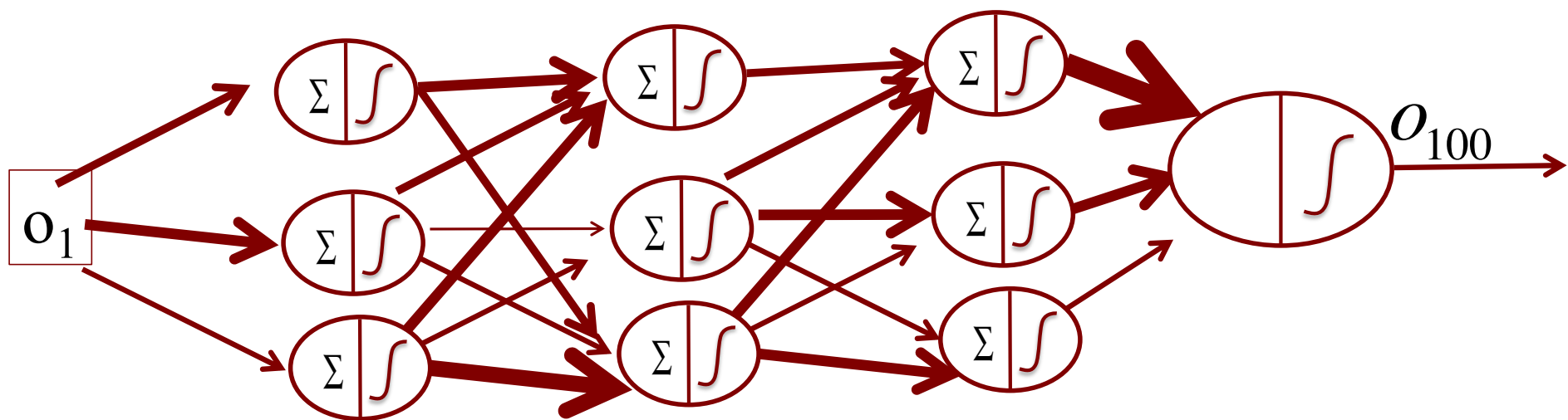
$$w_{a,b} \longleftarrow w_{a,b} - \alpha \frac{dE}{dw_{a,b}}$$

- $\alpha$  is between 0 and 1. Called the “learning rate”.
- Now we know how to propagate errors back through the network.
- Remember how to go forward?









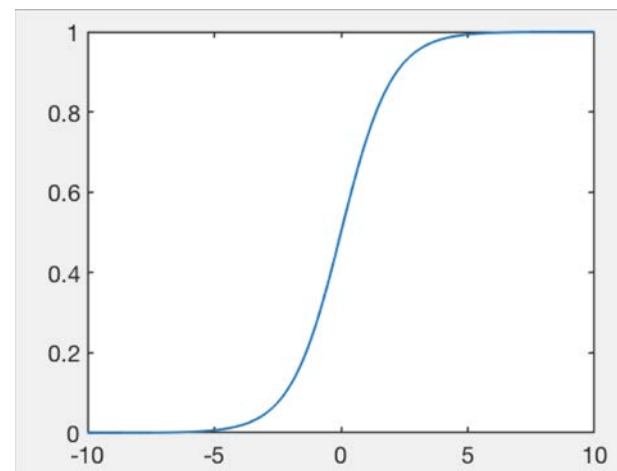


# Backpropagation

- Repeat going backwards (to calculate the gradients), adjusting the weights, and going forwards (to calculate the errors) over and over in order to learn.

# Convergence Problems

- NN's have **problems with convergence due to vanishing/exploding gradients and saddle points.**
- Vanishing gradients come from the flat part of the activation function.
- Exploding gradients happen when we realize that our gradient has vanished and so increase the learning rate and take huge step sizes to compensate (but then mess everything up!)
- Stick to  $10^{-5}$  to  $10^{-3}$  learning rate perhaps?



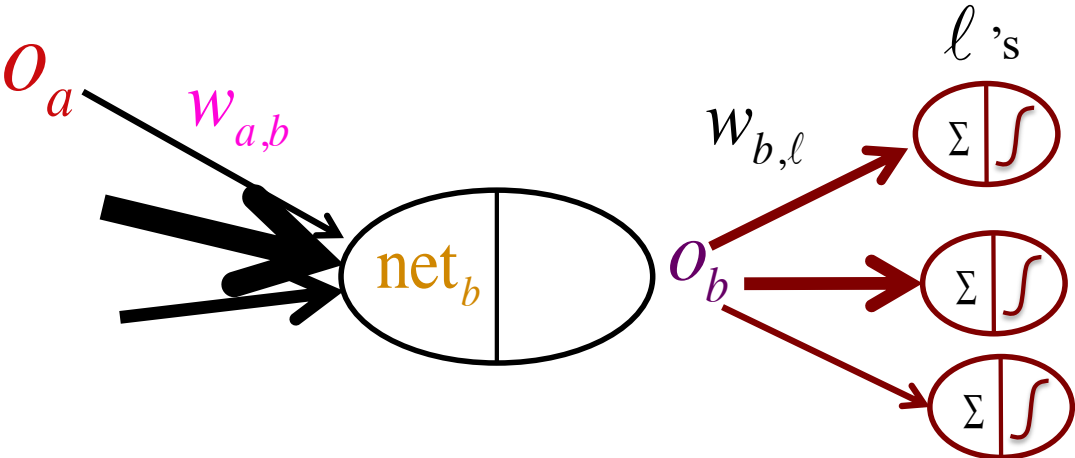
# Convergence Problems

- With the sigmoid activation, the derivatives of the input weights for each node are always **either all positive or all negative**. This is a limitation.

$$\frac{dE}{dw_{a,b}} = \frac{dE}{do_b} \frac{do_b}{d \text{net}_b} \frac{d \text{net}_b}{dw_{a,b}}$$

$$= \left( \sum_{\ell \in L} \delta_{\ell} w_{b,\ell} \right) o_b (1 - o_b) o_a$$

does not depend on a      positive, since all outputs of sigmoid are between 0 and 1.



# Convergence Problems

- Bottom line – most people do not use sigmoid-like activation functions, even though this is more biologically relevant.

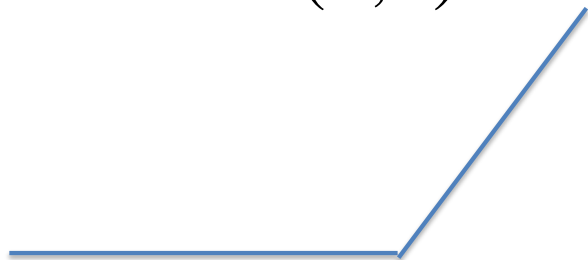
Sigmoid      $\sigma(x) = \frac{1}{1 + e^{-x}}$

Hyperbolic tangent      $\tanh(x)$

# Convergence Problems

## Rectified Linear Unit (ReLU)

$$\max(0, x)$$



Removes vanishing gradients when nodes are “activated,” meaning  $x > 0$ .

(Krizhevsky et al., 2012)

## Leaky ReLU

$$\max(0.1x, x)$$

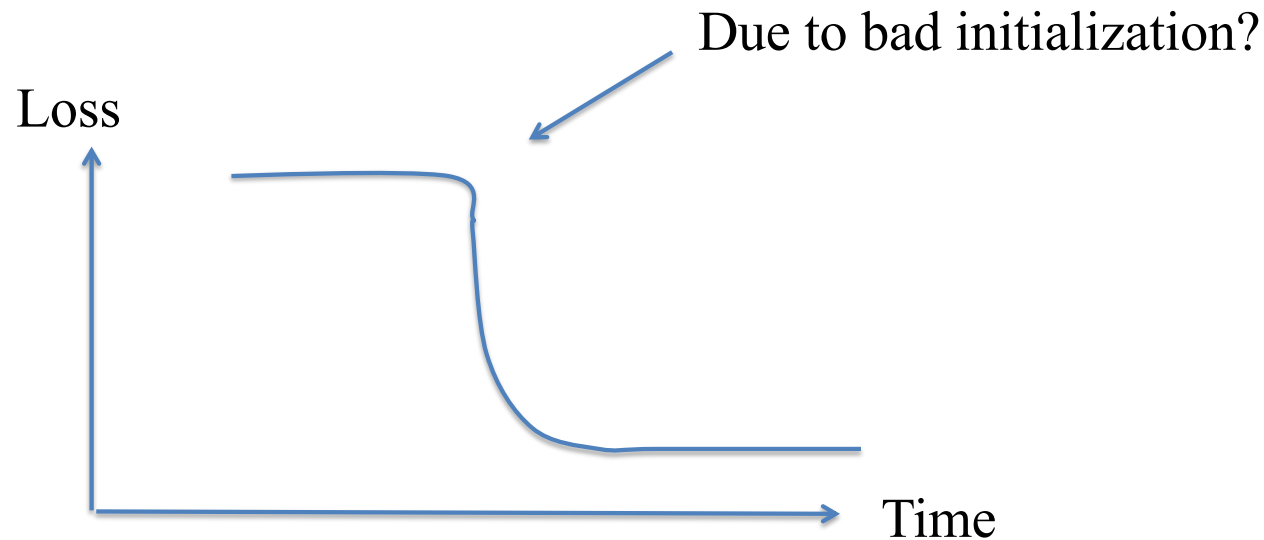


Removes vanishing gradients, but prefers that non-activated nodes be as “non-activated” as possible, which doesn’t make much sense.

(Mass et al., 2013; He et al., 2015)

# Convergence Problems

- **Initialization** of the networks weights is really important. I have no idea how to do it.



# Convergence Problems

- **Batch Normalization** (Ioffe and Szegedy, 2015) is a step in a NN that:
  - Normalizes the outputs  $o_i$  of several nodes (a “mini-batch”) in the same layer (as usual, subtract the mean of the  $o_i$ ’s divide by their standard deviation).
  - Includes the mean and standard deviation as separate parameters to be learned.
  - Usually the normalization is before the nonlinear activation function.
  - This adds regularization and helps gradient flow in the network but sometimes it messes things up.

# Data Augmentation



Chinese Lantern Festival, Cary NC, 2017



# Data Augmentation

- include artificial data, such as horizontal flips, rotations, resized, cropped training images, change contrast and brightness, distortion, etc.



Chinese Lantern Festival, Cary NC, 2017

# Convolutional NN's

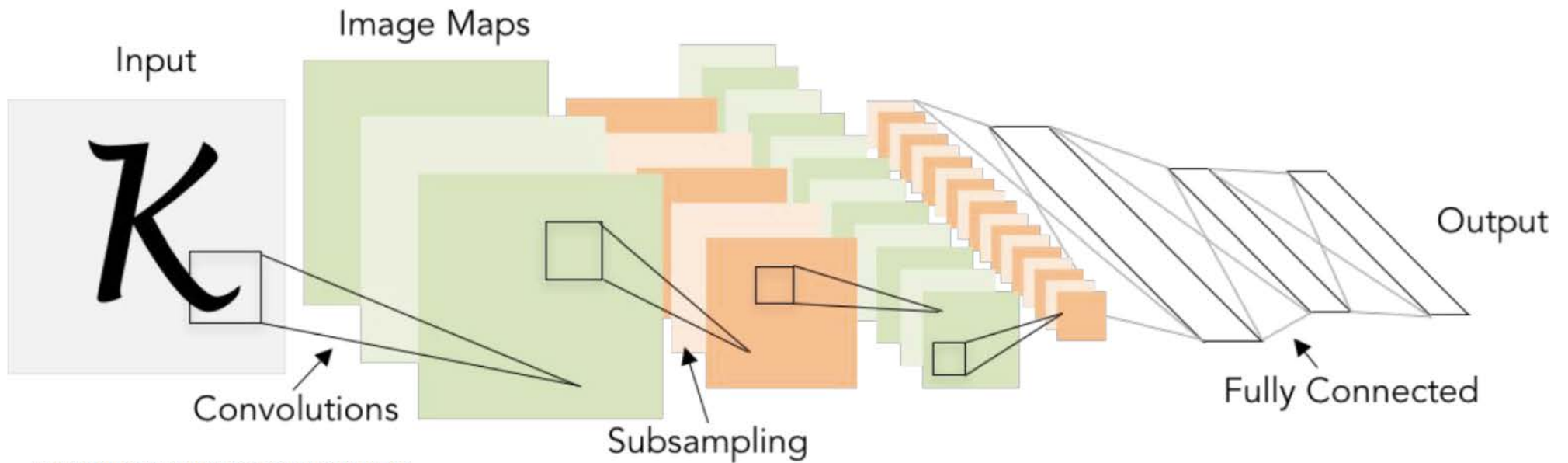
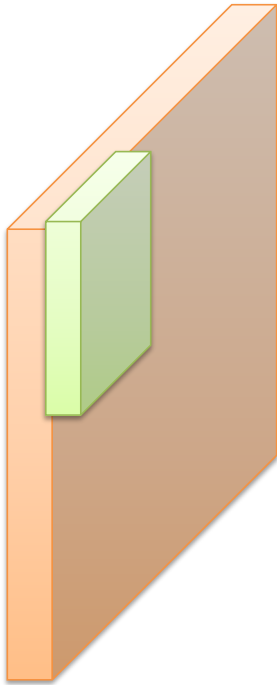


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

Image from LeCun et al 1998, reproduced also from Li, Johnson, Yeung, 2017

# Convolutional NN's

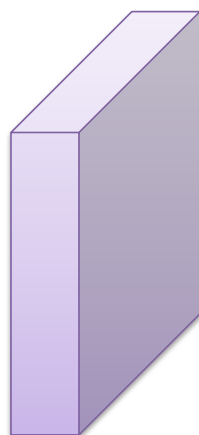
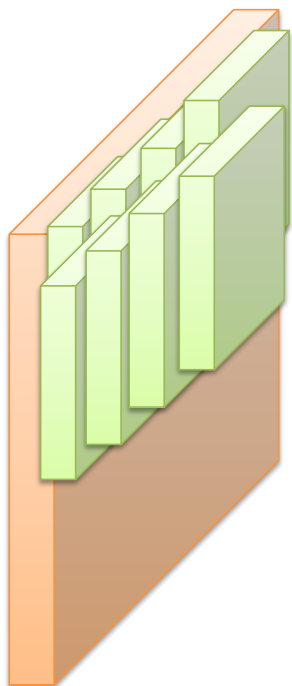
- Convolve means to slide the filter over all spatial locations, and sum up the filter weights times the input.



- An edge filter will detect edges.

# Convolutional NN's

- Convolve means to slide the filter over all spatial locations, and sum up the filter weights times the input.

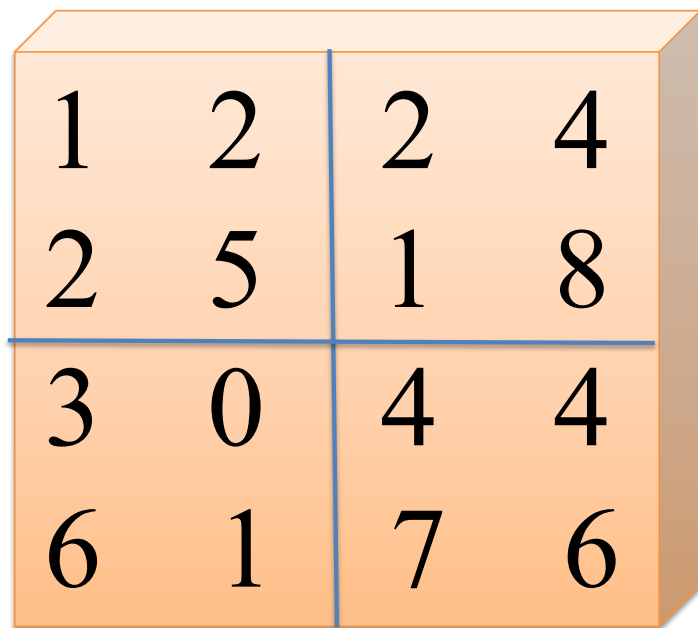


- Stride of 5 means we step by 5's when we convolve.

The following layer is smaller by a factor of 5.

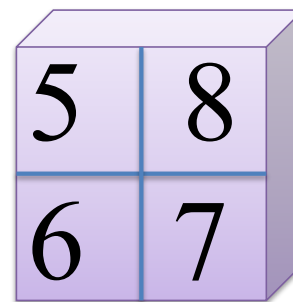
# Convolutional NN's

- **Max pooling** means to convolve with a max function.
- Intuitively keeps track of whether an earlier filter has detected something.



1	2	2	4
2	5	1	8
3	0	4	4
6	1	7	6

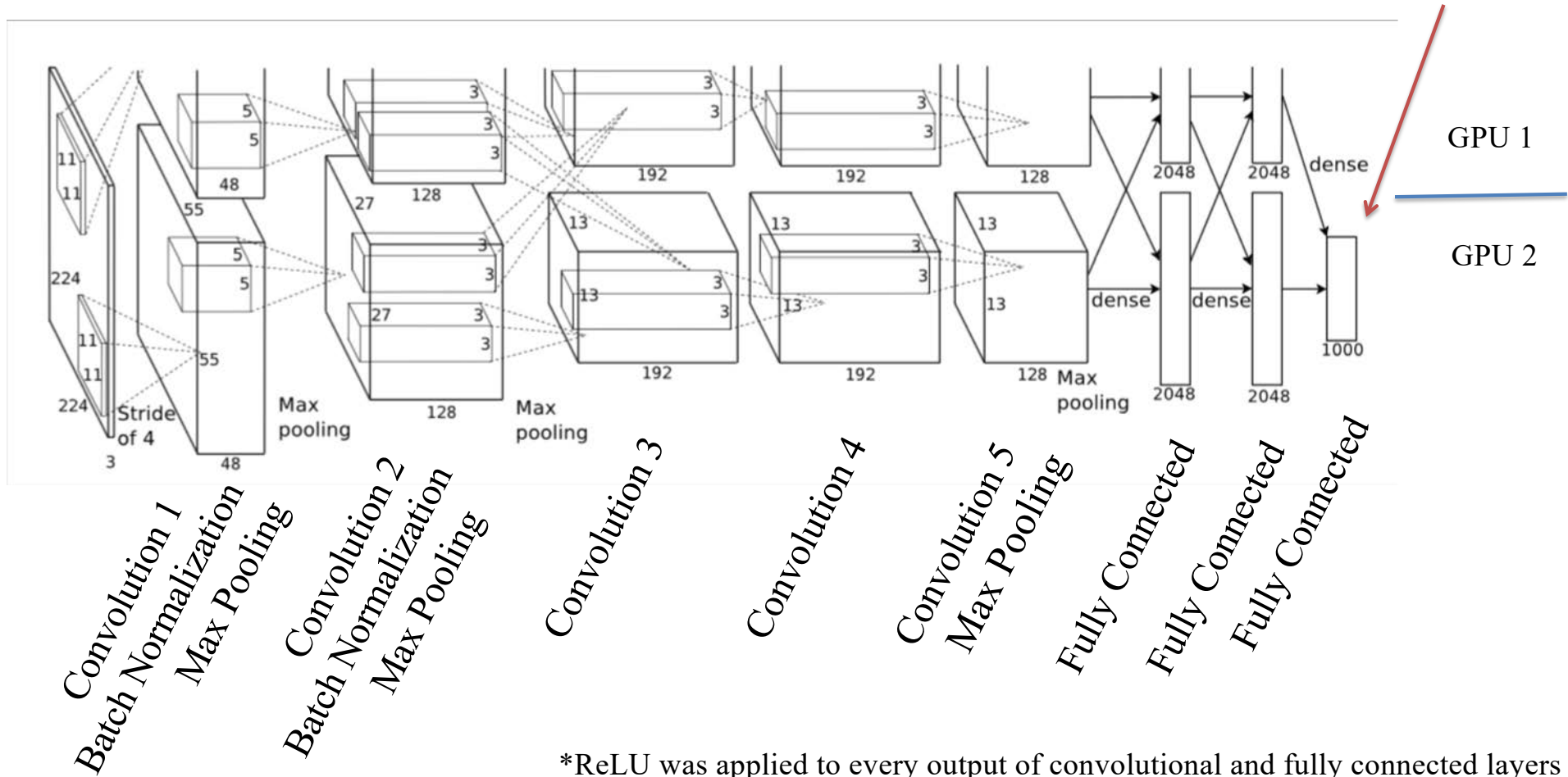
2 x 2 max pool filter and stride 2



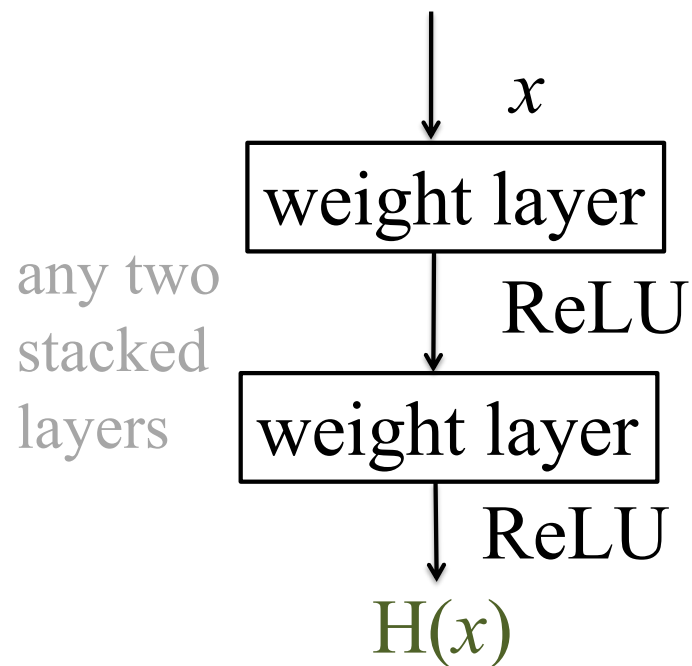
5	8
6	7

- AlexNet (Krizhevsky et al. 2012)

softmax over 1000 classes



# Residual Nets (He et al., 2016)

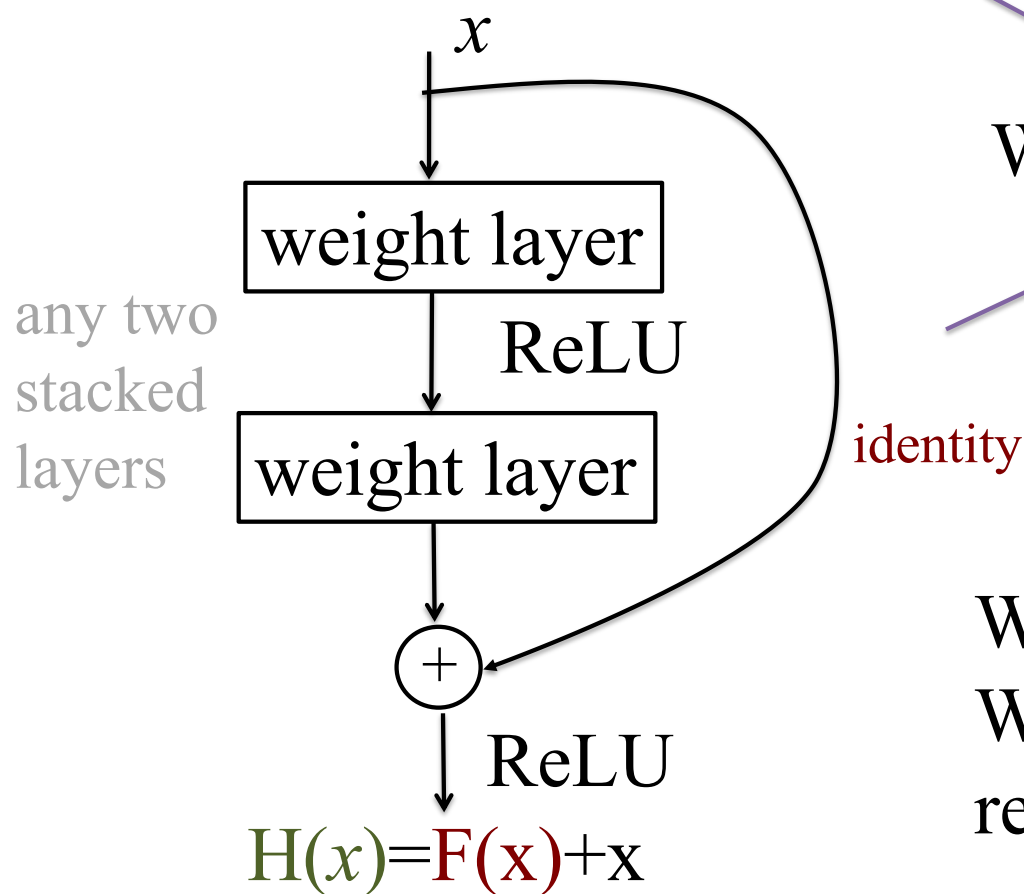


~~We hope to fit  $H(x)$ .~~

Slides recreated from Kaiming He's tutorial

[http://kaiminghe.com/icml16tutorial/icml2016\\_tutorial\\_deep\\_residual\\_networks\\_kaiminghe.pdf](http://kaiminghe.com/icml16tutorial/icml2016_tutorial_deep_residual_networks_kaiminghe.pdf)

# Residual Nets



~~We hope to fit  $H(x)$ .~~

We hope to fit  $F(x)$ .  
We are now learning a residual of identity.



# Residual Nets

- By adding  $x$ , the derivative of the error with respect to  $x$  increases by 1. Thus, less vanishing derivatives.
- Allowed networks to go much deeper than before.  
“From 10 to 1000 layers”

$$H(x) = F(x) + x$$

# Dropout (Srivastava et al., JMLR 2014)

- In each forward pass, for each neuron, with probability  $p$ , set all of its output weights to 0.
- $p$  is a hyperparameter, usually  $p = 0.5$ .
- During testing, use all nodes.

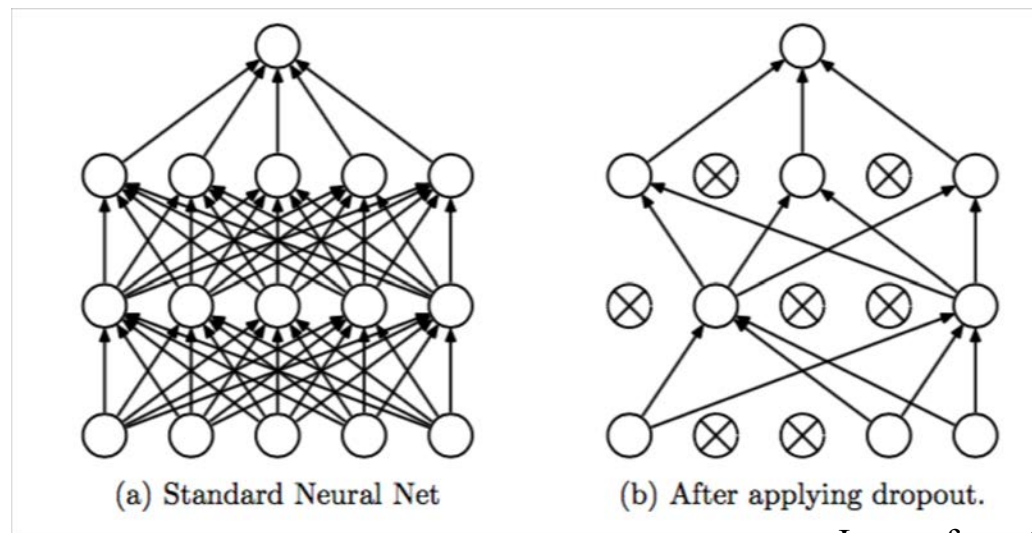


Image from Srivastava et al JMLR 2014)

# Dropout (Srivastava et al., JMLR 2014)

- As if we are training exponentially many “sub” models. Similar idea to bagging (averaging many separately trained models together).
- Creates a redundant encoding.

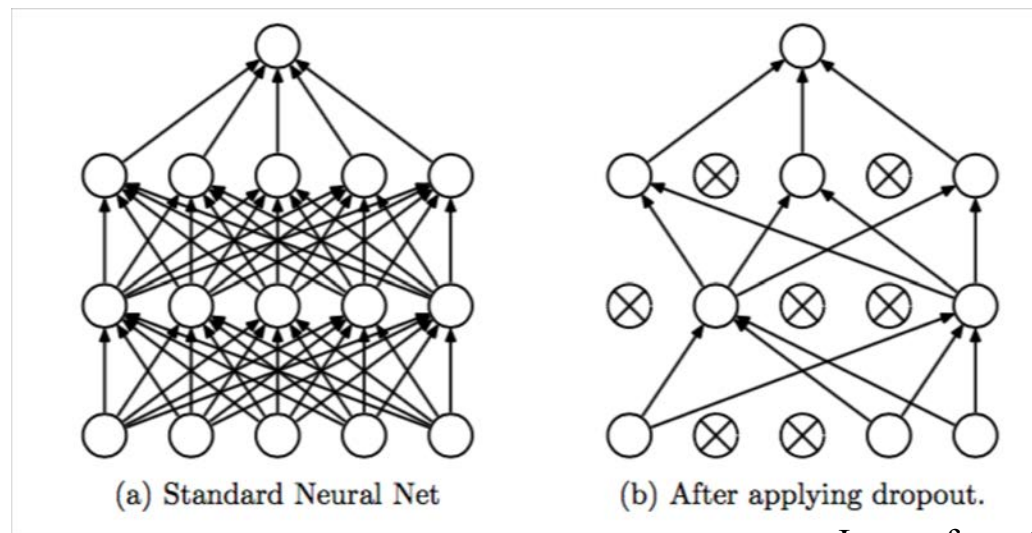
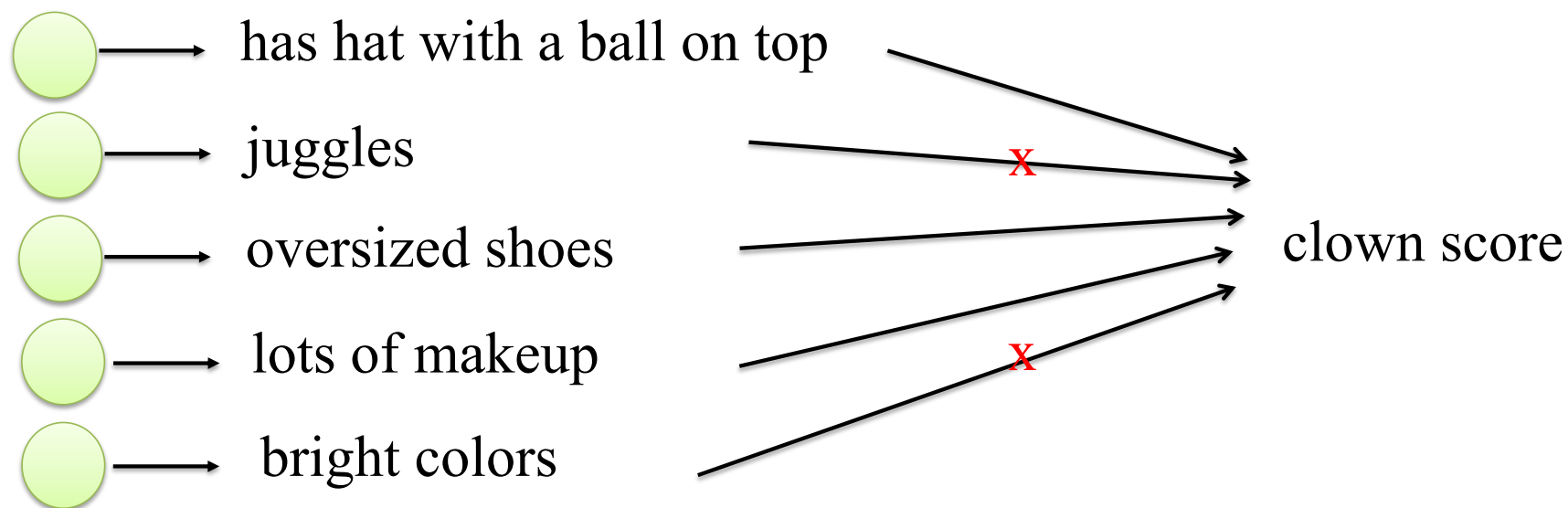


Image from Srivastava et al JMLR 2014)

# Dropout (Srivastava et al., JMLR 2014)

- As if we are training exponentially many “sub” models. Similar idea to bagging (averaging many separately trained models together).
- Creates a redundant encoding.

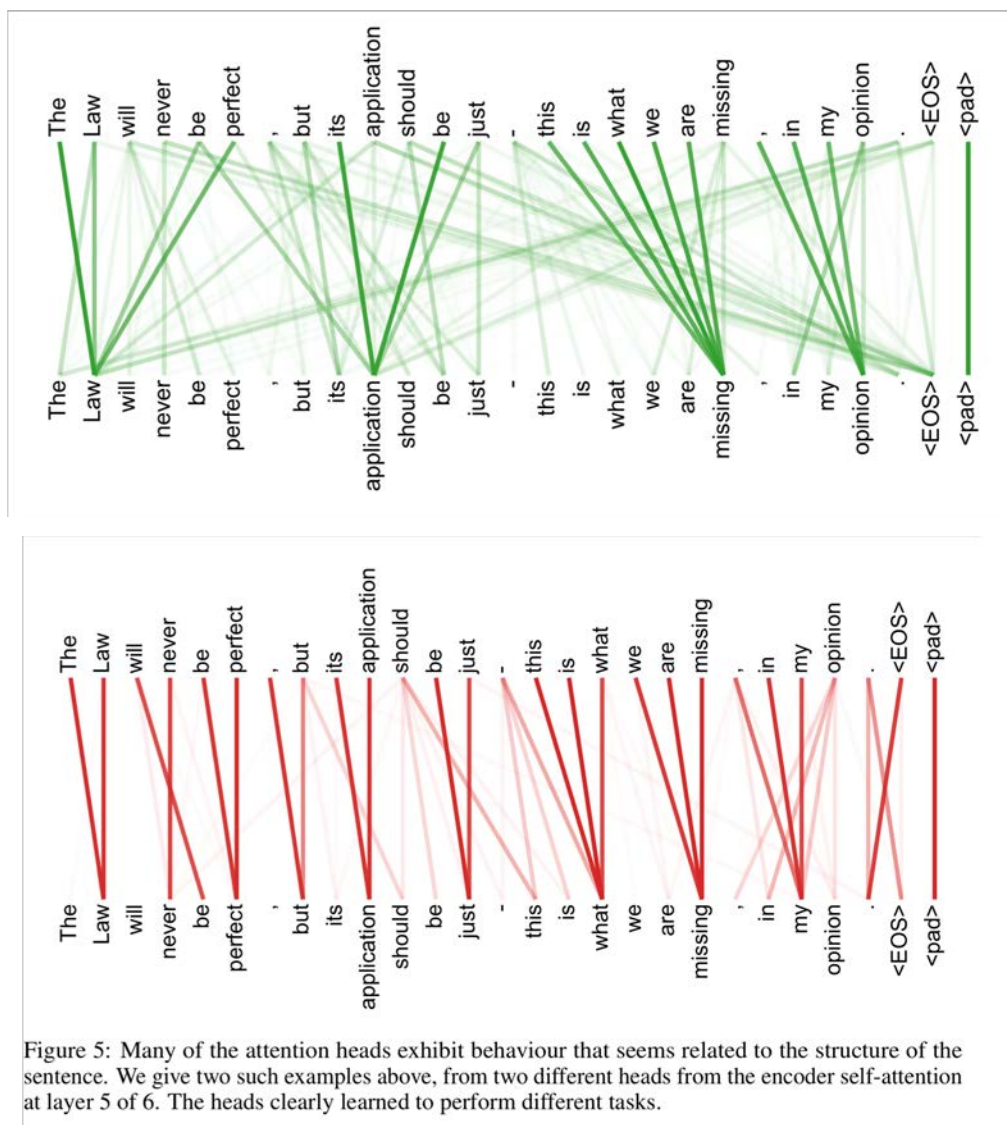


# A Big Bag of Tricks

- Dropout
- Batch Normalization
- Data Augmentation
- Residual Networks
- Convolutional Networks
- Activation Functions (ReLU, Leaky ReLU)
- Initialization
- Max Pooling
- :

# RNN's have recently been surpassed by Transformers

- RNN's used to be the standard for machine translation (sequence-to-sequence). These were useful because they kept a memory of longer-term relationships between words.
- The new standard is the transformer, which places several “multi-head attention nodes” together. The transformer compares all words to all other words in the sentence. It is not recursive.



“Attention is all you need” Vaswani et al 2017

# GANs – Generative Adversarial Networks

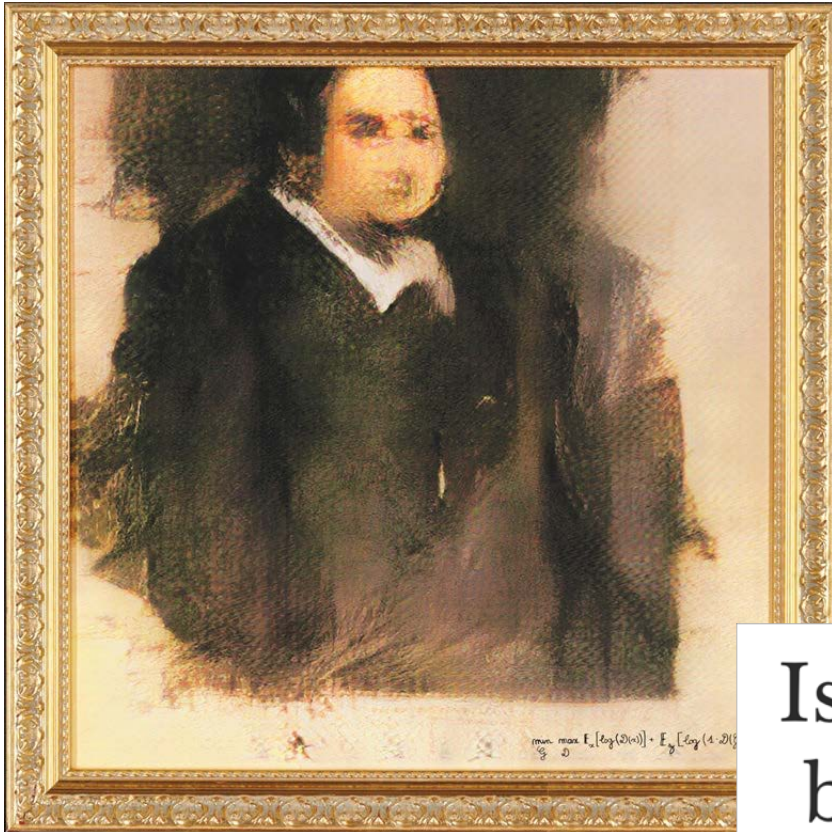
From Goodfellow et al 2014:

In other words,  $D$  and  $G$  play the following two-player minimax game with value function  $V(G, D)$ :

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (1)$$

- GANs are actor-critic models
- They produce realistic-looking images/data (some ethical/legal concerns)
- Used commonly for AI artwork





## Is artificial intelligence set to become art's next medium?

AI artwork sells for \$432,500 — nearly 45 times its high estimate — as Christie's becomes the first auction house to offer a work of art created by an algorithm

A



B



C



D



Figure adapted from L. Gatys et al. "[A Neural Algorithm of Artistic Style](#)" (2015) by Google AI Blog

# Neural networks

- Advantages:
  - highly expressive nonlinear models
  - have advances in computer vision and speech that other methods have not achieved
  - can capture latent structure within the hidden layers
- Disadvantages
  - can get stuck in local optima, could produce bad solutions
  - black box
  - lots of tuning parameters (e.g., the structure of the network)
- There is a recent bag of tricks

Thanks