

# Statistical Natural Language Processing [COMP0087]

*Introduction to neural networks  
and backpropagation*

**Vasileios Lampos**

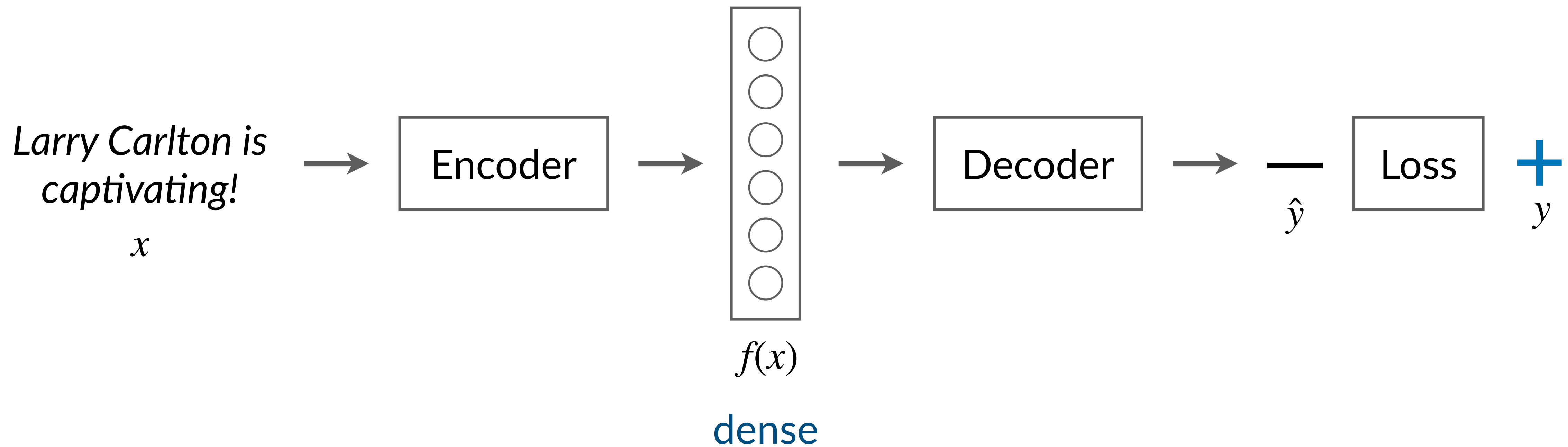
Computer Science, UCL



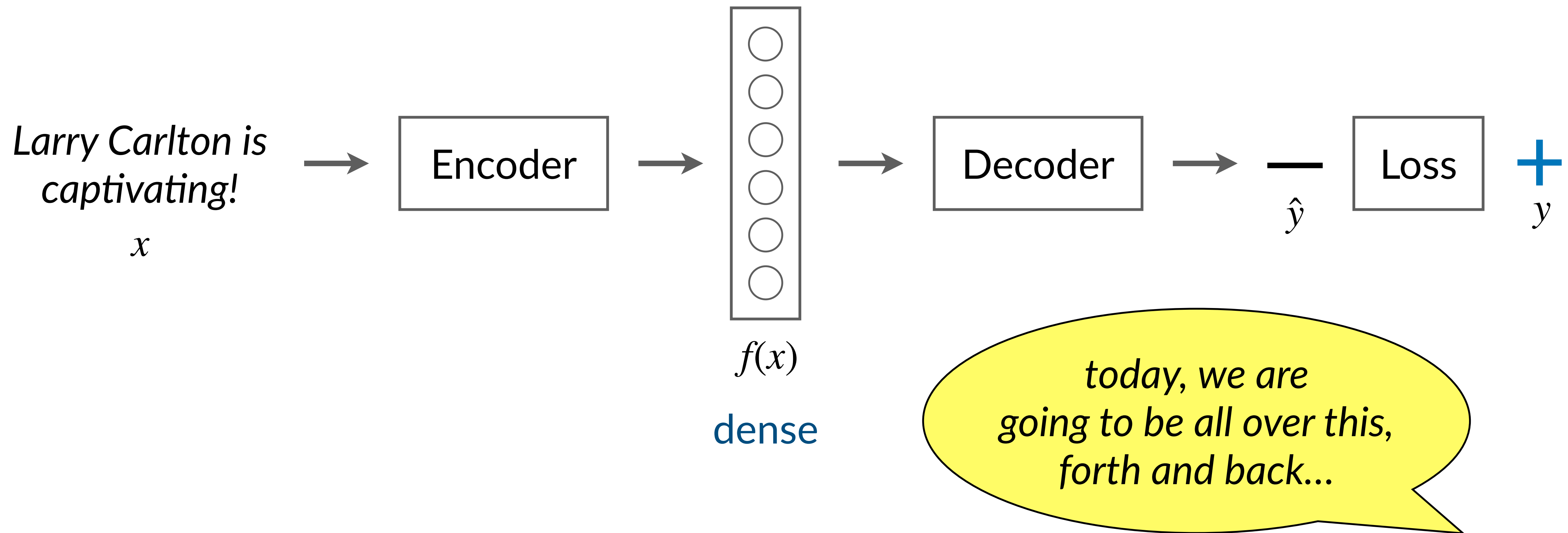
# About this lecture

- ▶ In this lecture:
  - Introductory neural network concepts
  - Inference and training (*backpropagation*) with feedforward neural networks
- ▶ Reading / Lecture partly based on: Chapter 7 of “*Speech and Language Processing*” (SLP) by Jurafsky and Martin (2023) — [web.stanford.edu/~jurafsky/slp3/](http://web.stanford.edu/~jurafsky/slp3/)
- ▶ For those of you who want to have the slides in front of them during the lecture, there is a clipped / early version at [lampos.net/teaching](http://lampos.net/teaching) (*non clipped / slightly refined version will be added after the lecture*)

# The NLP view (for this lecture)

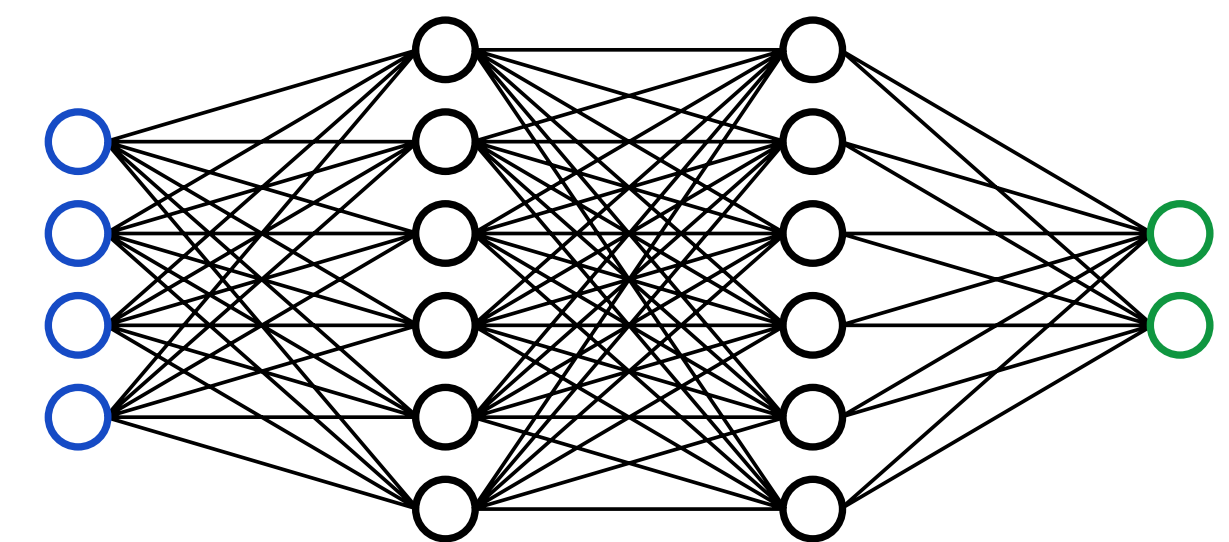
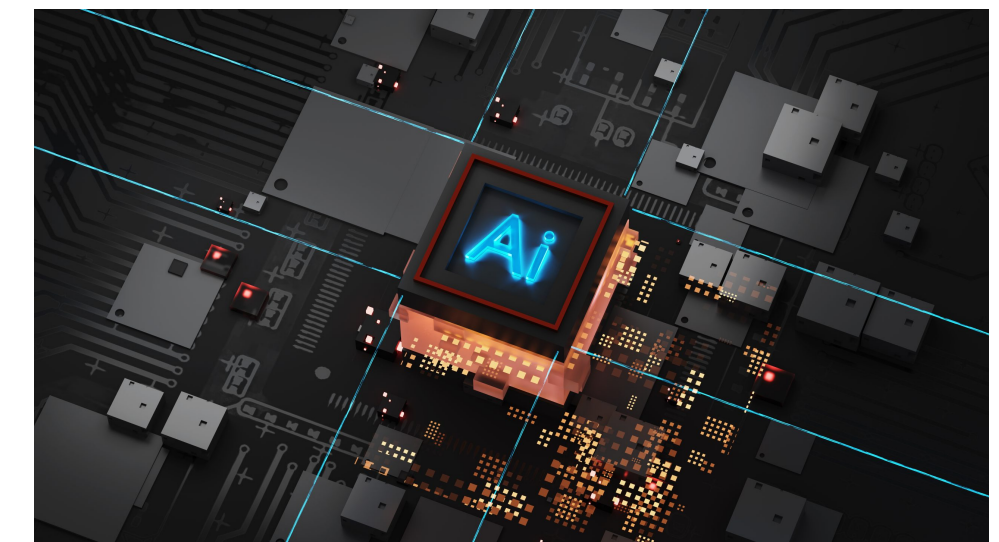
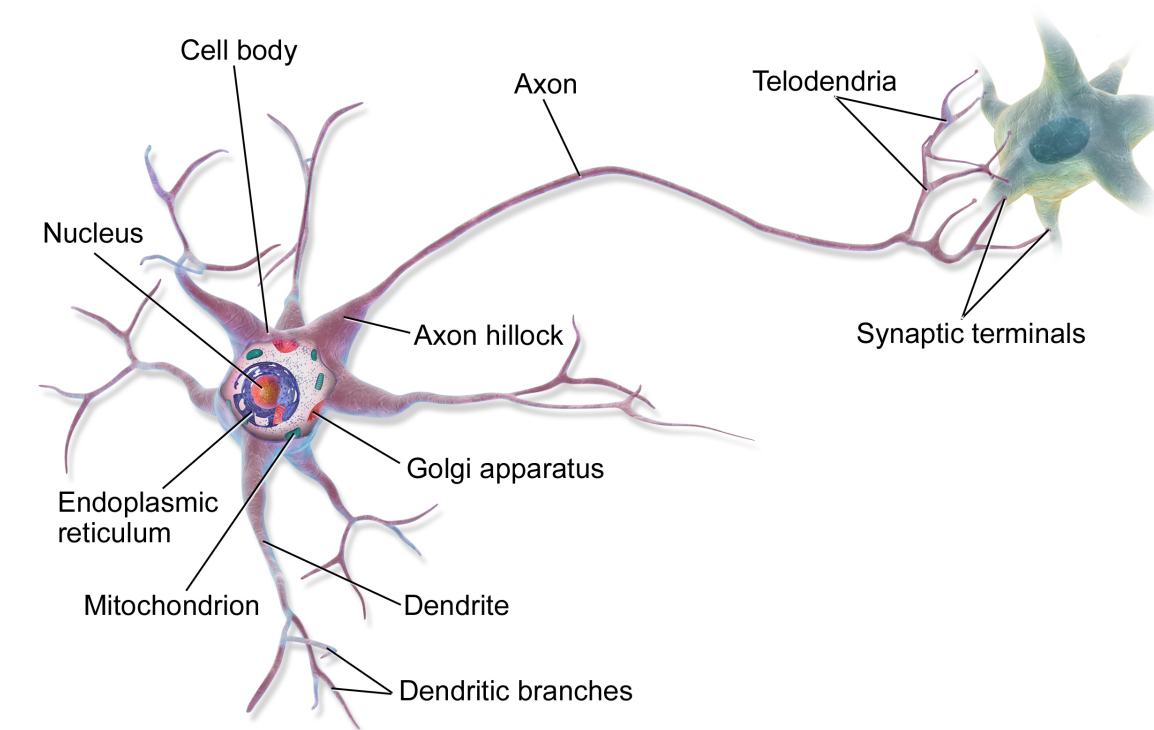


# The NLP view (for this lecture)



# Artificial neural networks – A few introductory remarks

- ▶ Artificial Neural Networks (NNs)  $\neq$  biological neural networks *until we actually obtain a complete understanding about how the human brain operates!*
- ▶ NNs are powerful learning functions / universal approximators, e.g. standard multi-layer feedforward networks with as few as one hidden layer are capable of approximating any (*Borel measurable*) function – and we are aware of this for almost 40 years (Hornik, Stinchcombe and White, 1989, [doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8))
- ▶ NB: Good understanding of logistic regression? Easy to understand today's lecture and fundamentals about NNs in a few seconds. *Otherwise it might take a few minutes.*





# Background task – Sentiment classification

Sentiment?

*Wow, I love the sound of this acoustic guitar!*

→ + (positive)

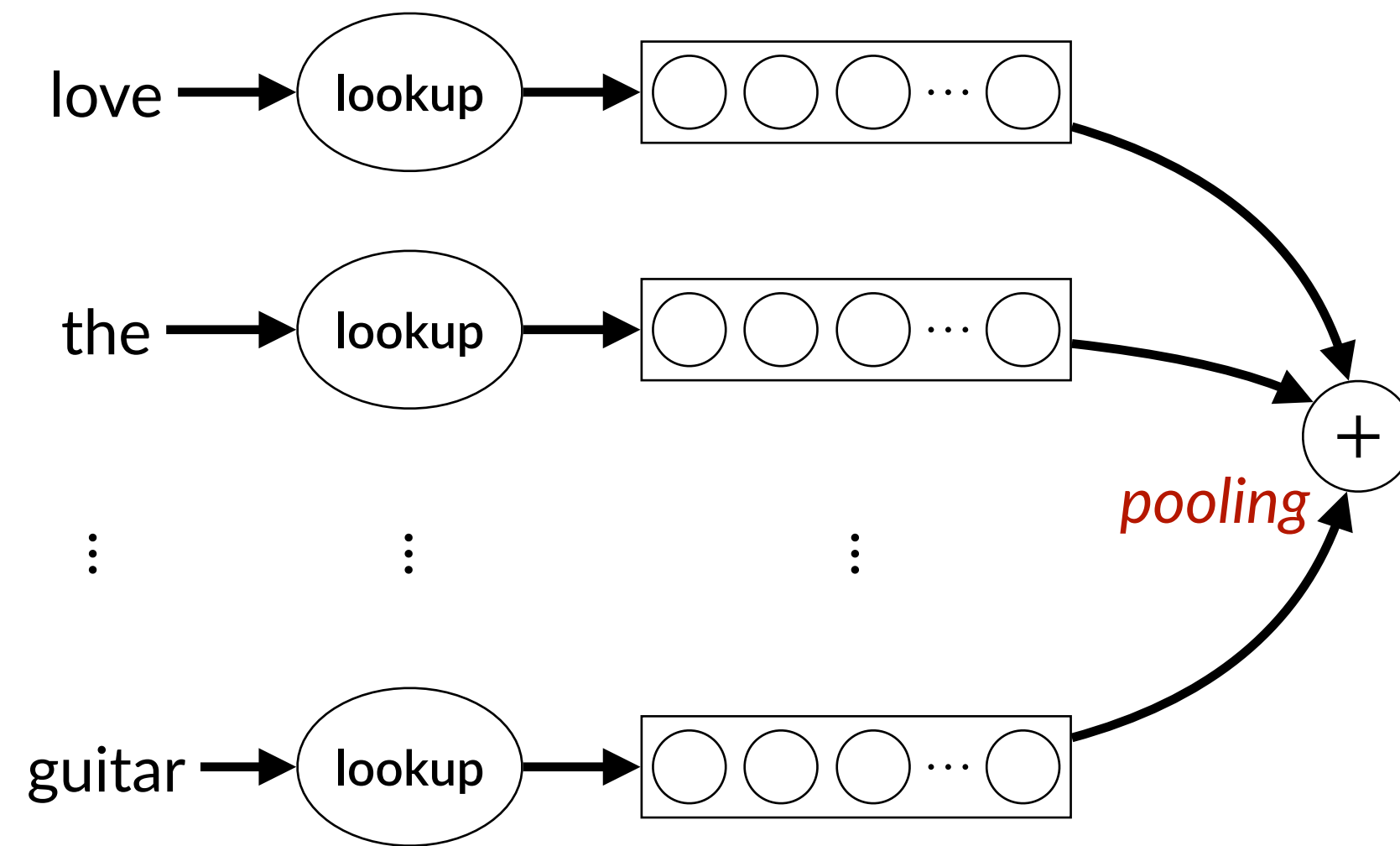
*It was just another uneventful Marvel movie!*

→ - (negative)

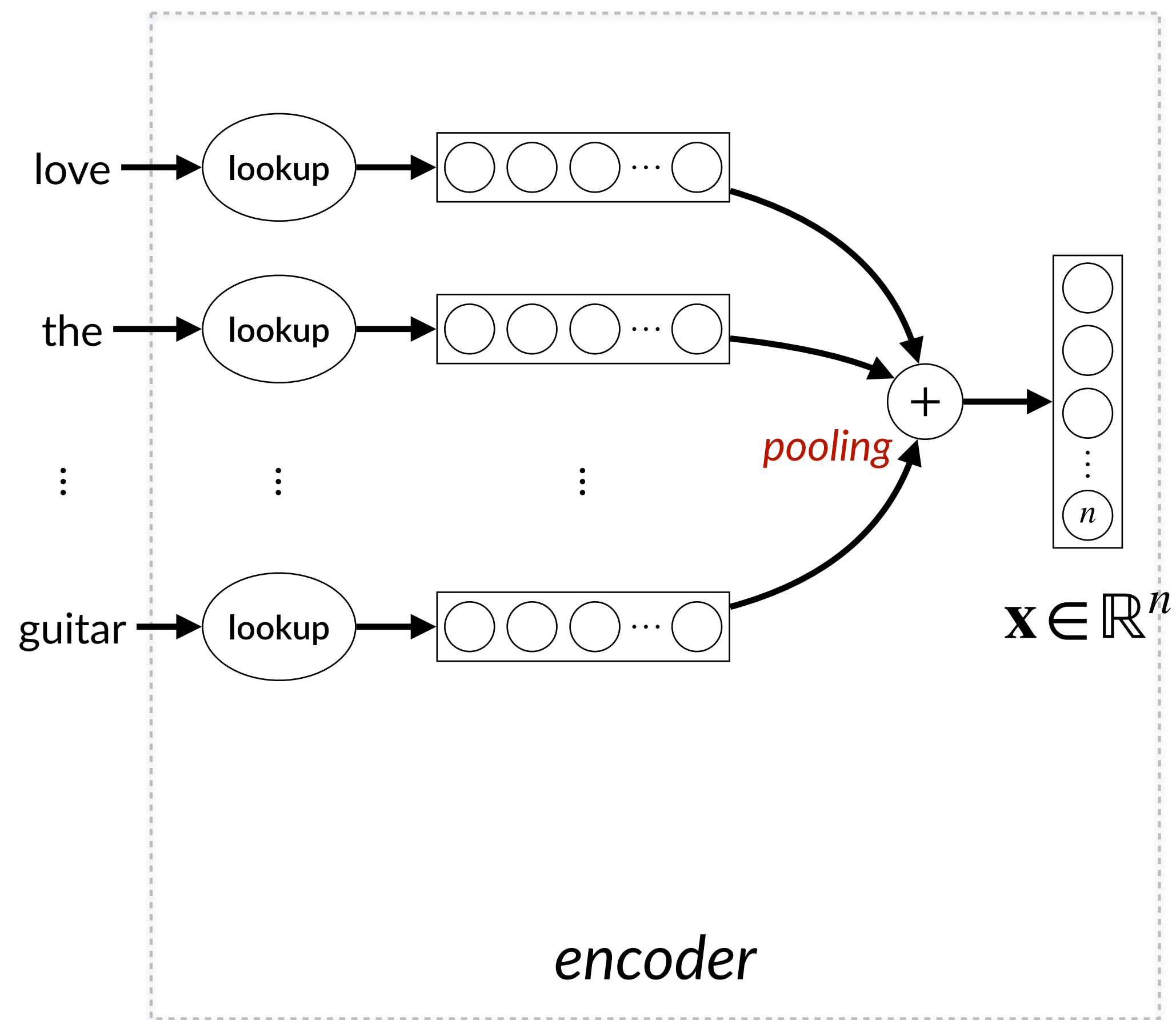
*Can't say I loved this performance, but I didn't dislike it either.*

→ neutral

# Neural network – A simplified encoder–decoder architecture

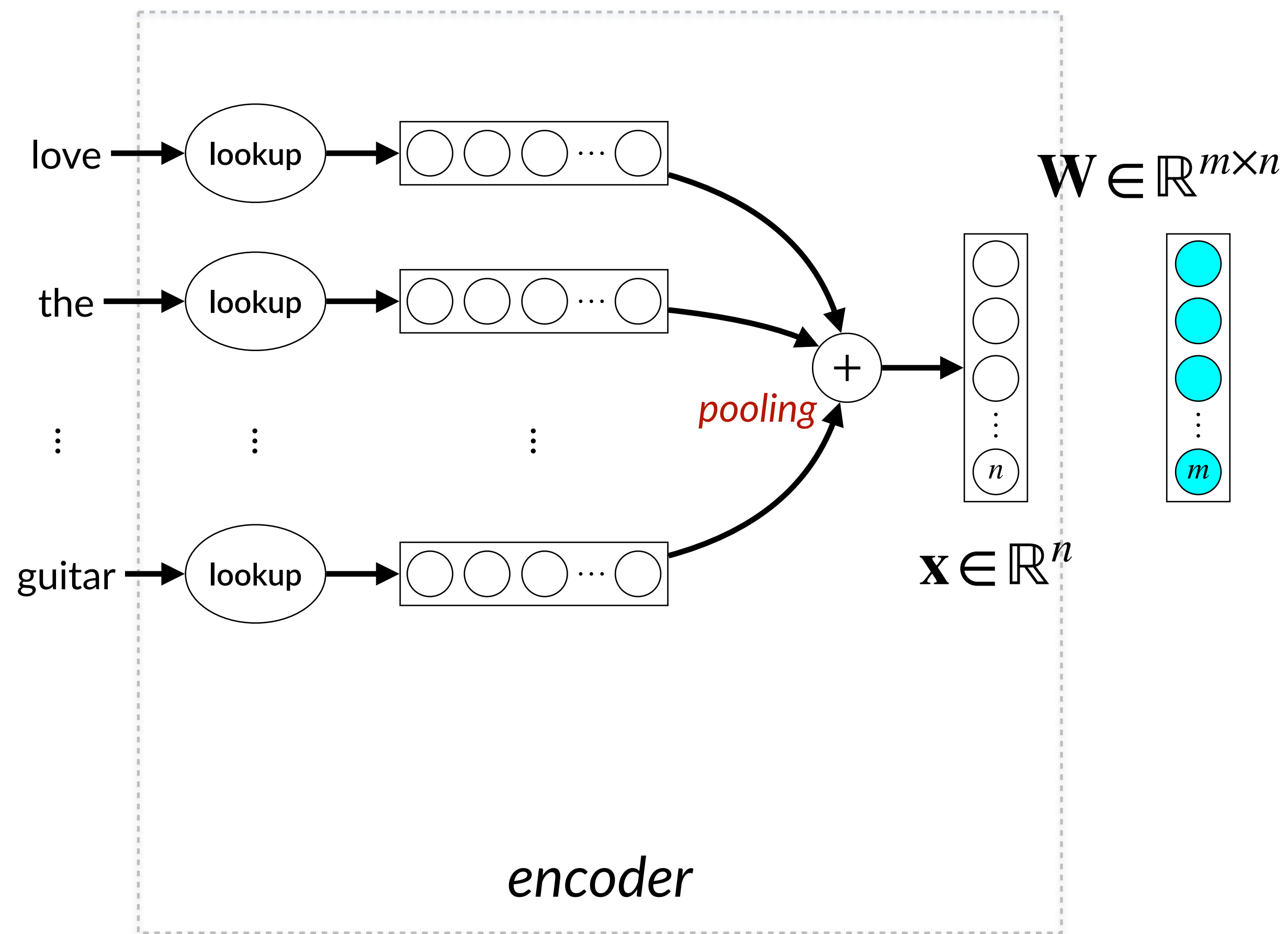


# Neural network – A simplified encoder–decoder architecture

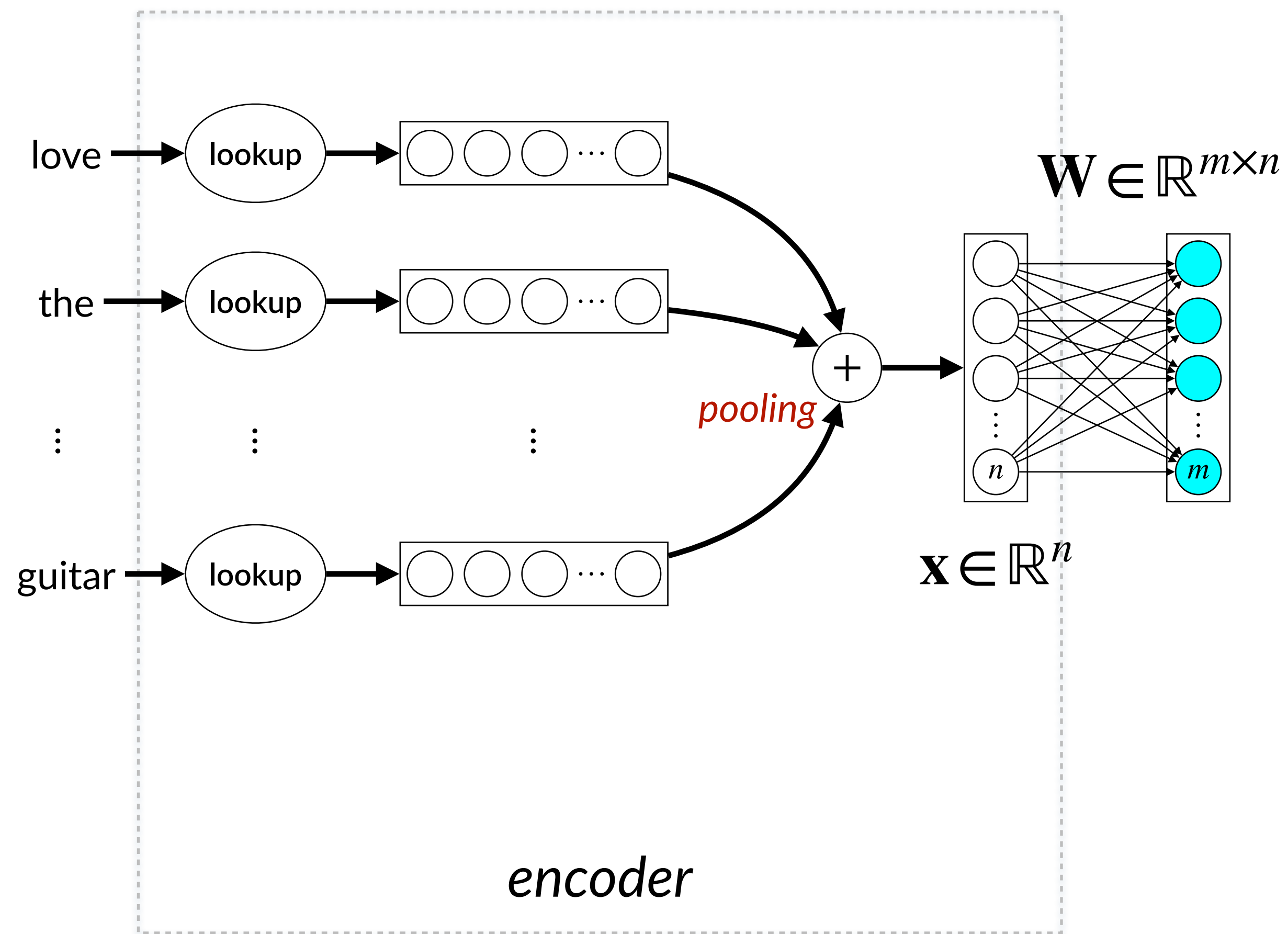




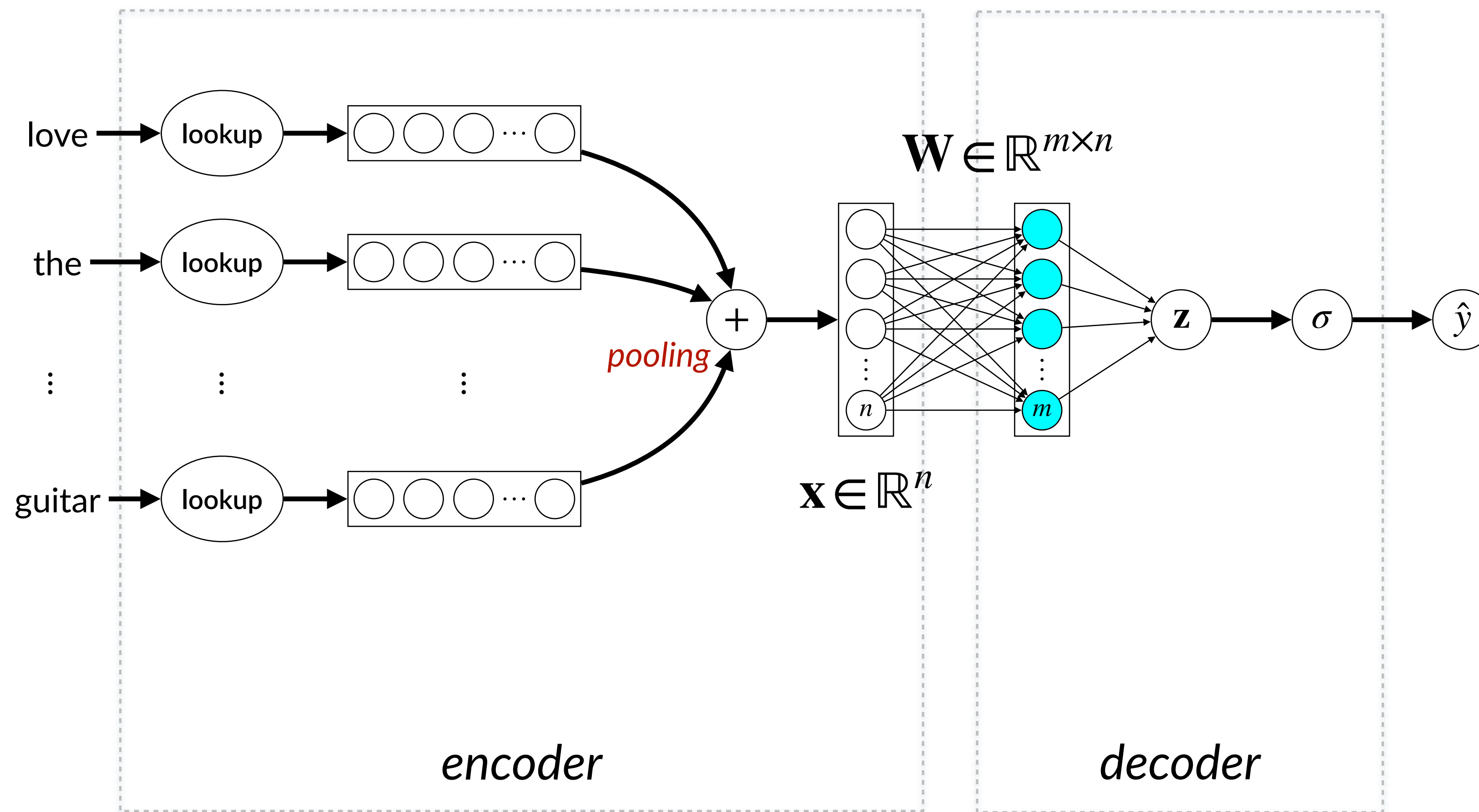
# Neural network – A simplified encoder–decoder architecture



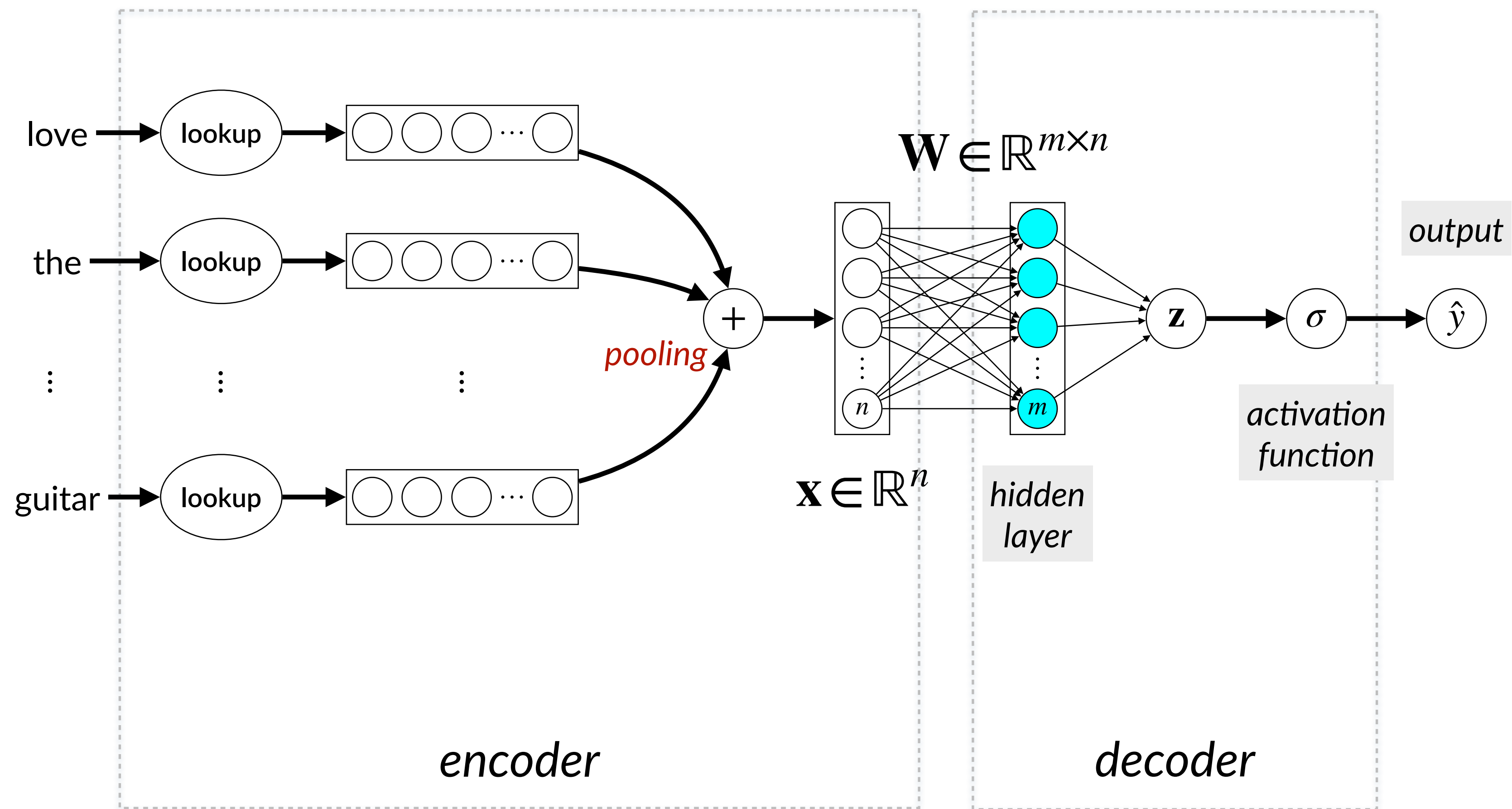
# Neural network – A simplified encoder–decoder architecture



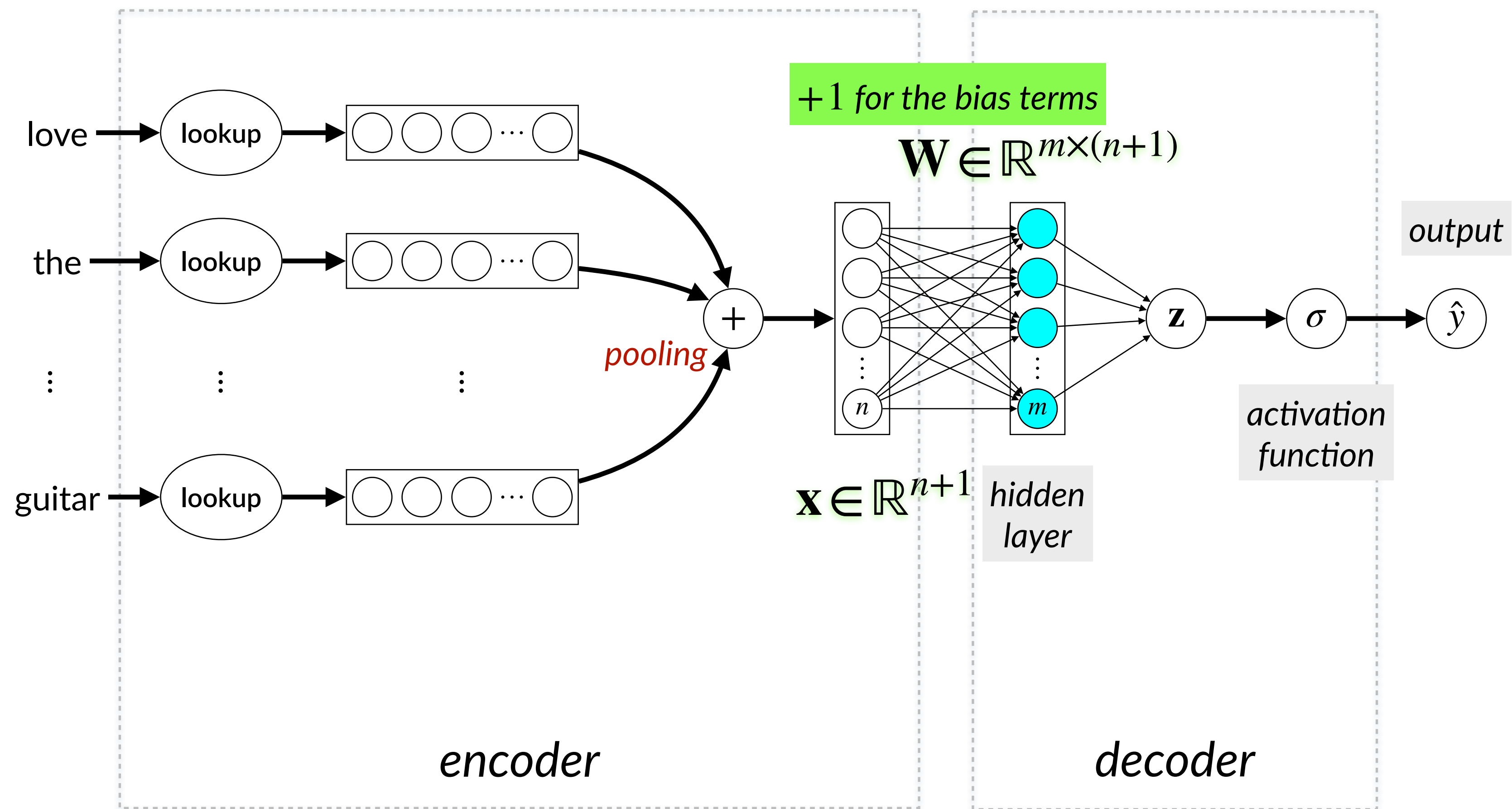
# Neural network – A simplified encoder–decoder architecture



# Neural network – A simplified encoder–decoder architecture

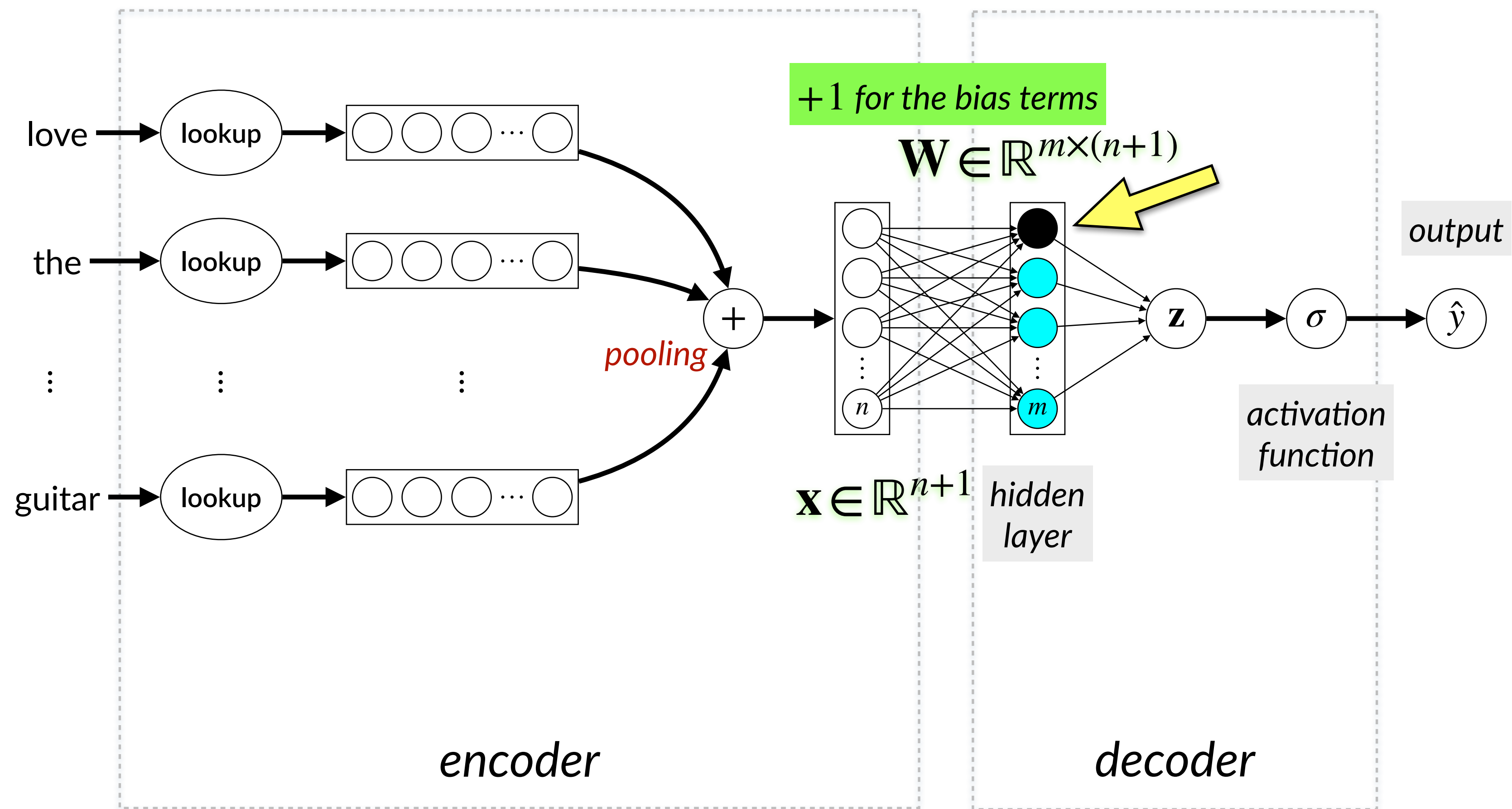


# Neural network – A simplified encoder–decoder architecture



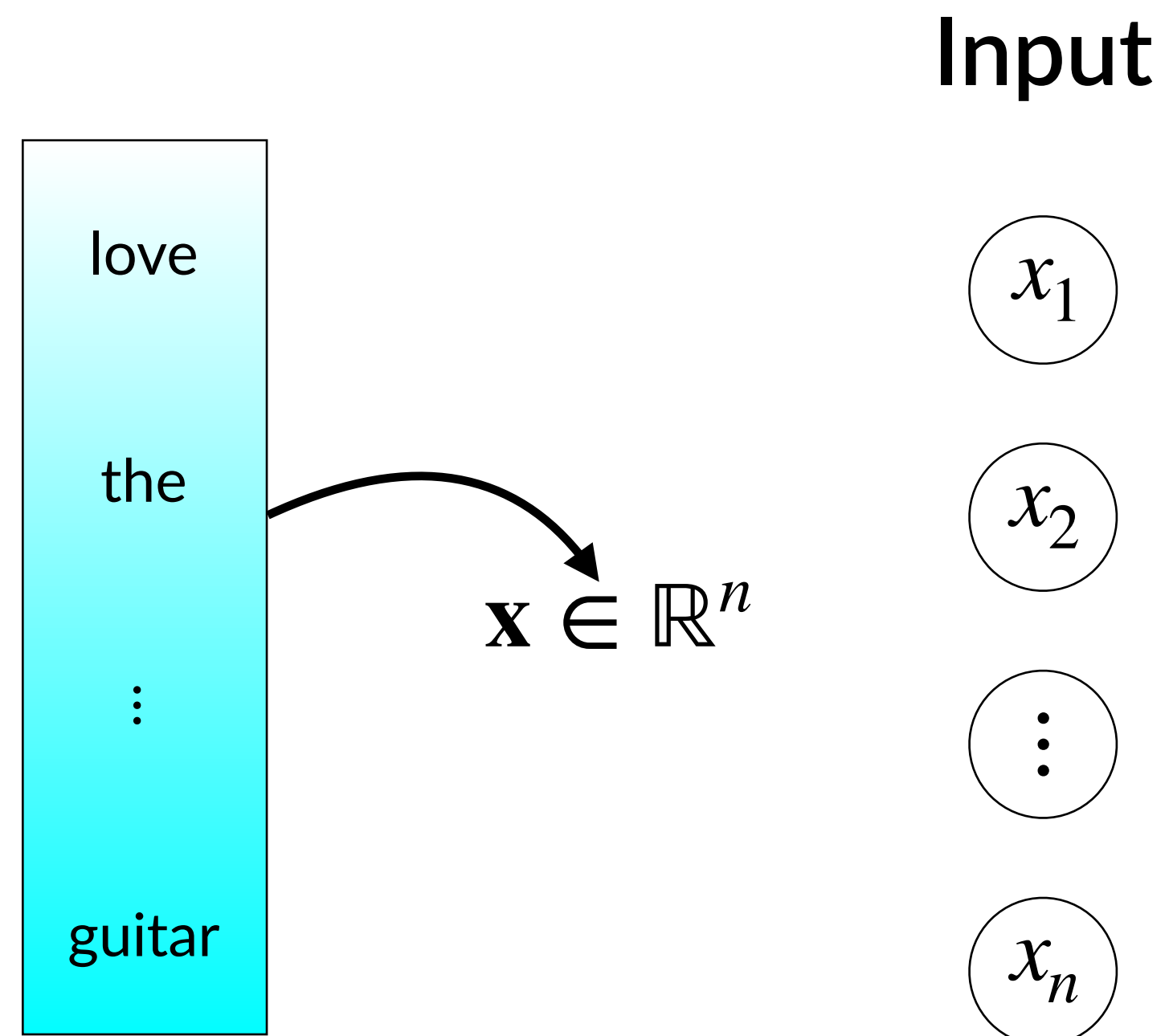


# Neural network – A simplified encoder–decoder architecture



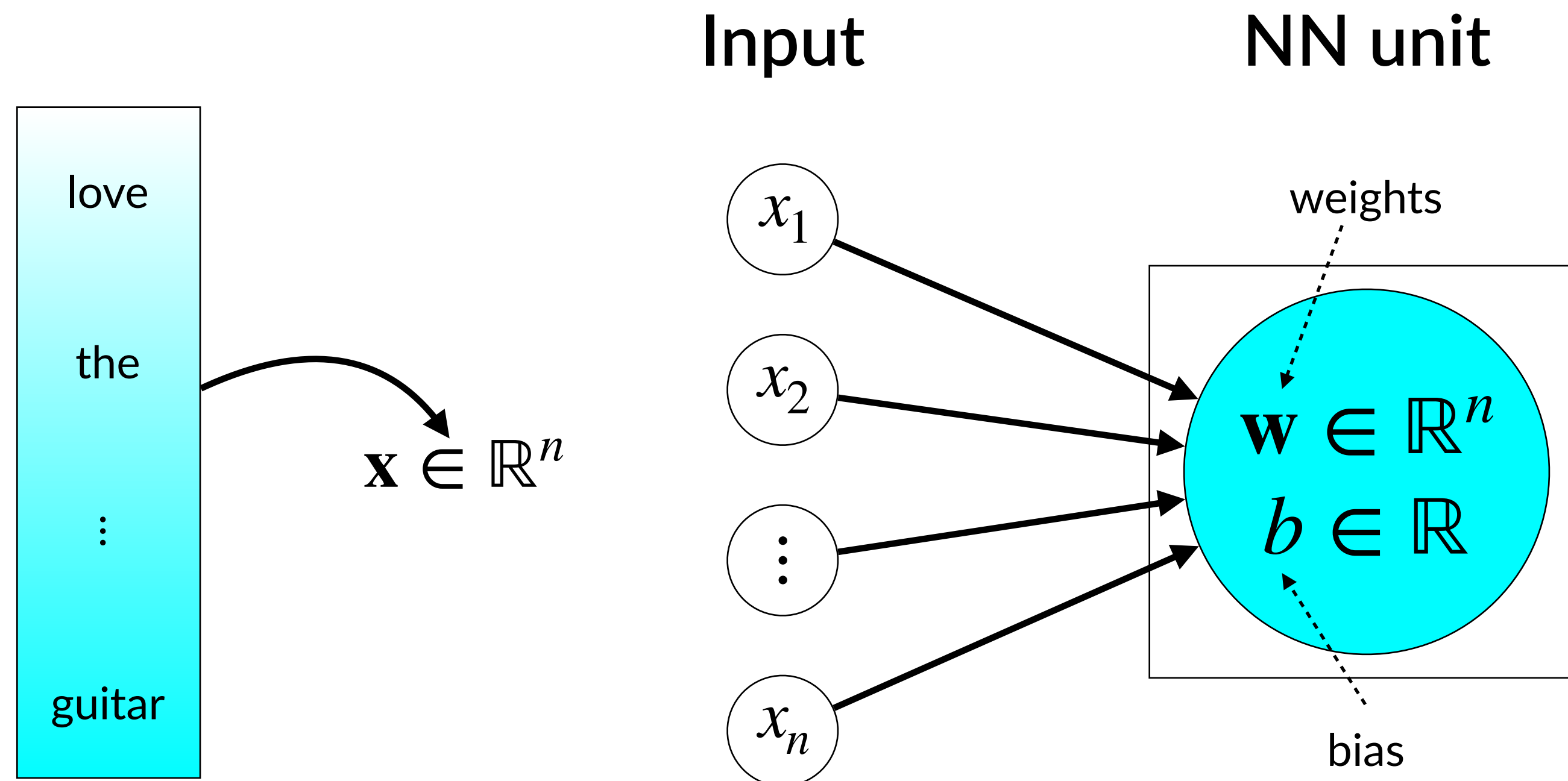


# One neural computation unit (of a hidden layer)



$$\mathbf{x} = [x_1 \quad \dots \quad x_n]$$

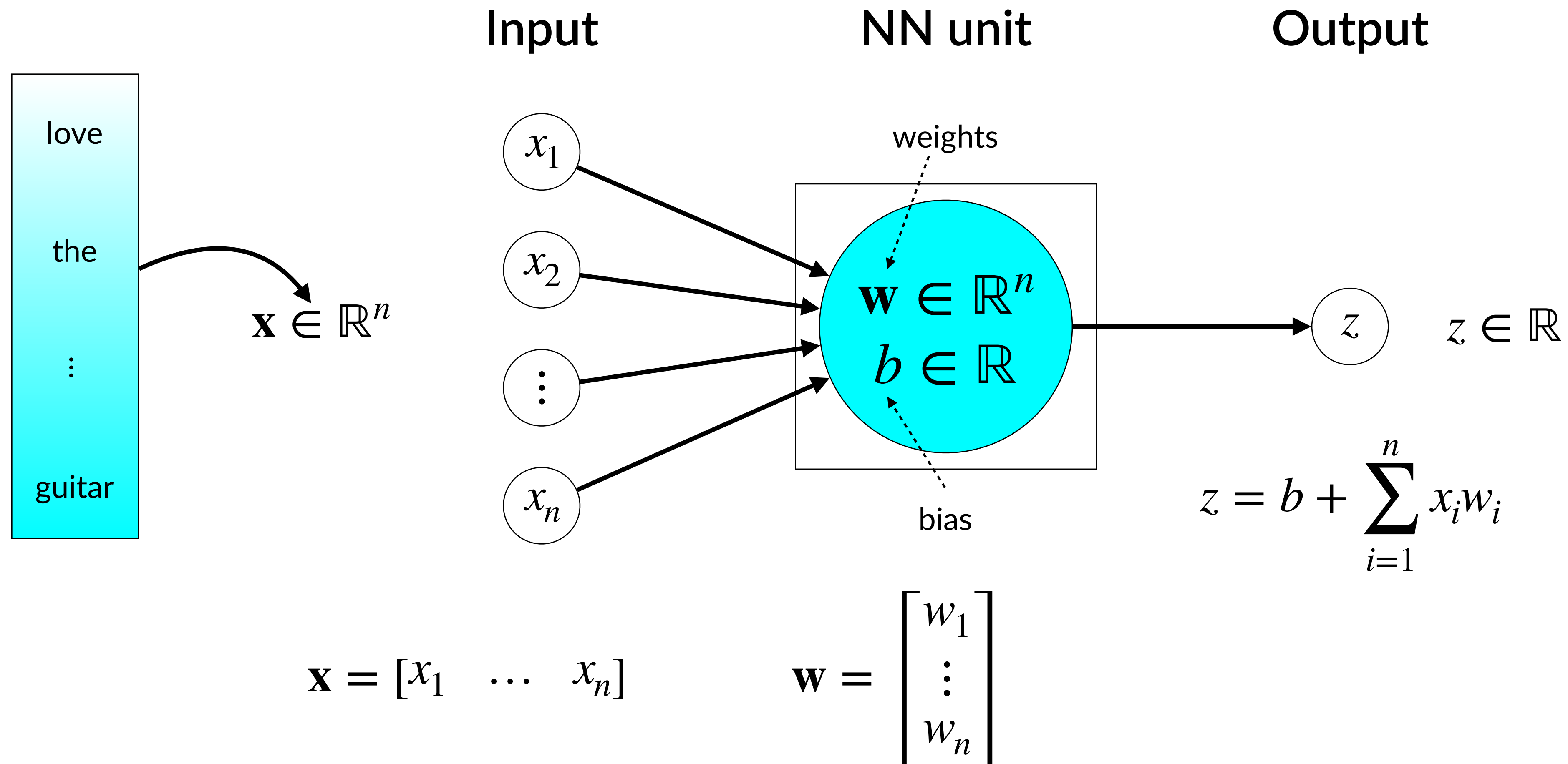
# One neural computation unit (of a hidden layer)



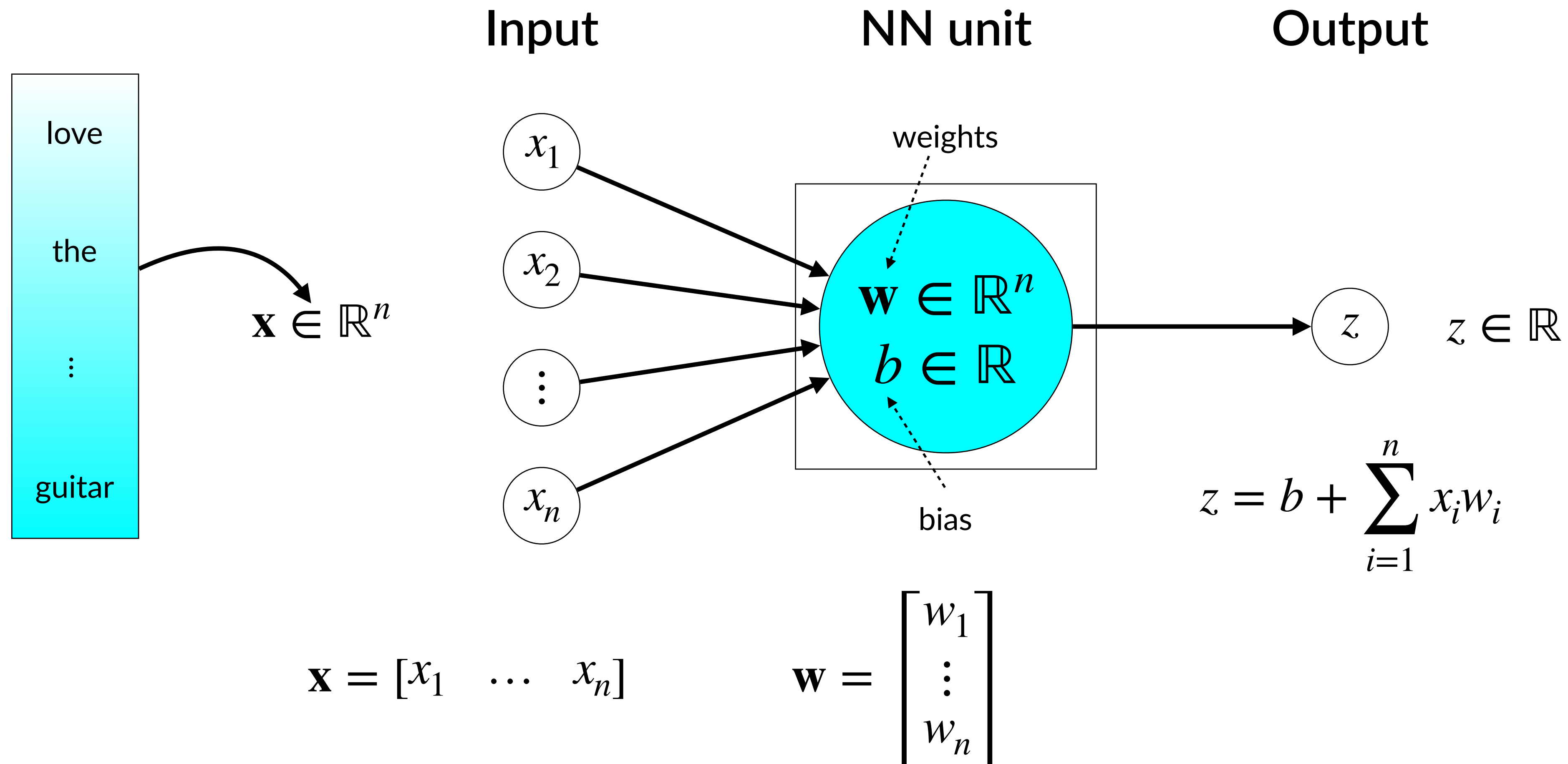
$$\mathbf{x} = [x_1 \quad \dots \quad x_n]$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

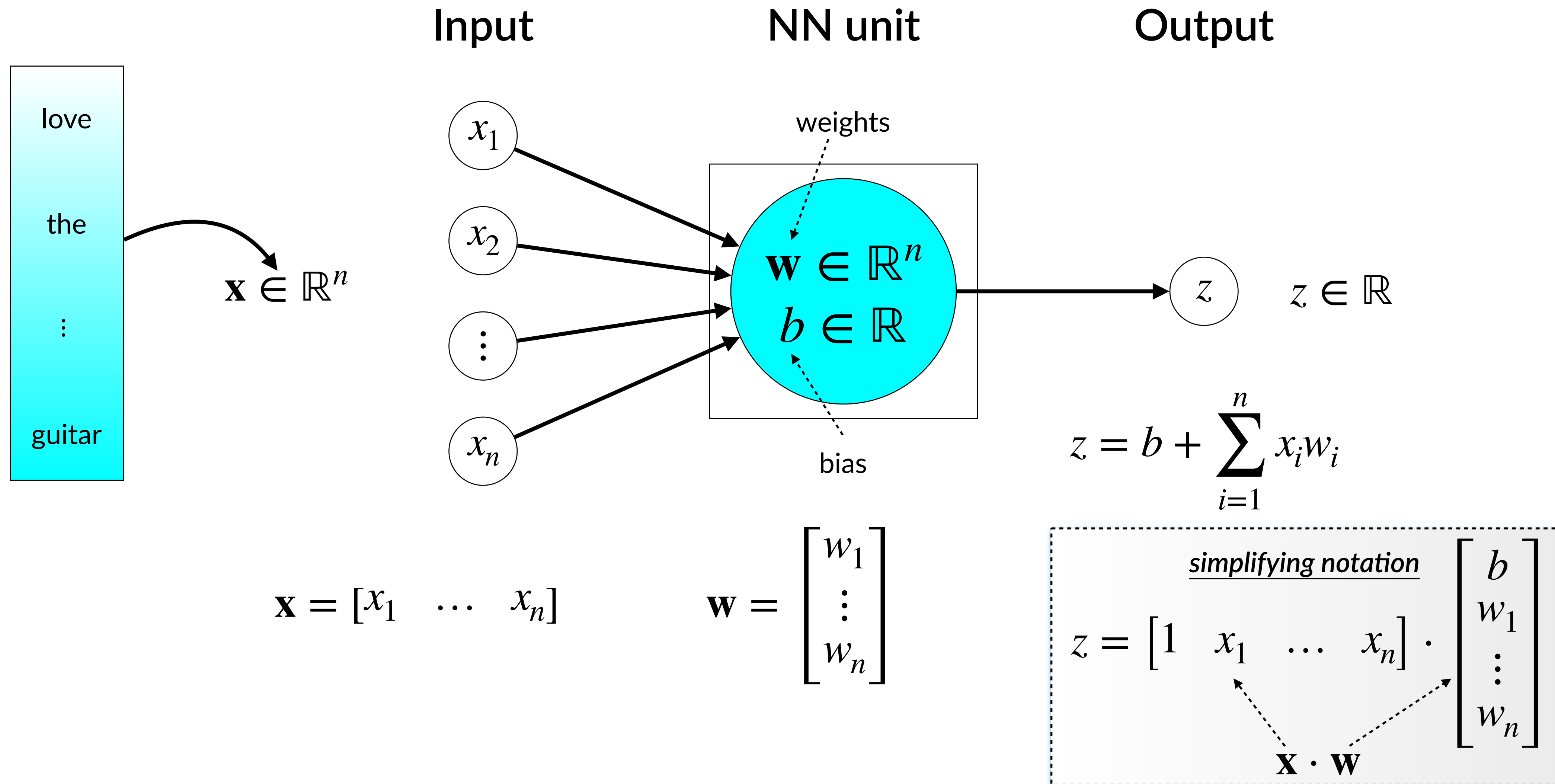
# One neural computation unit (of a hidden layer)



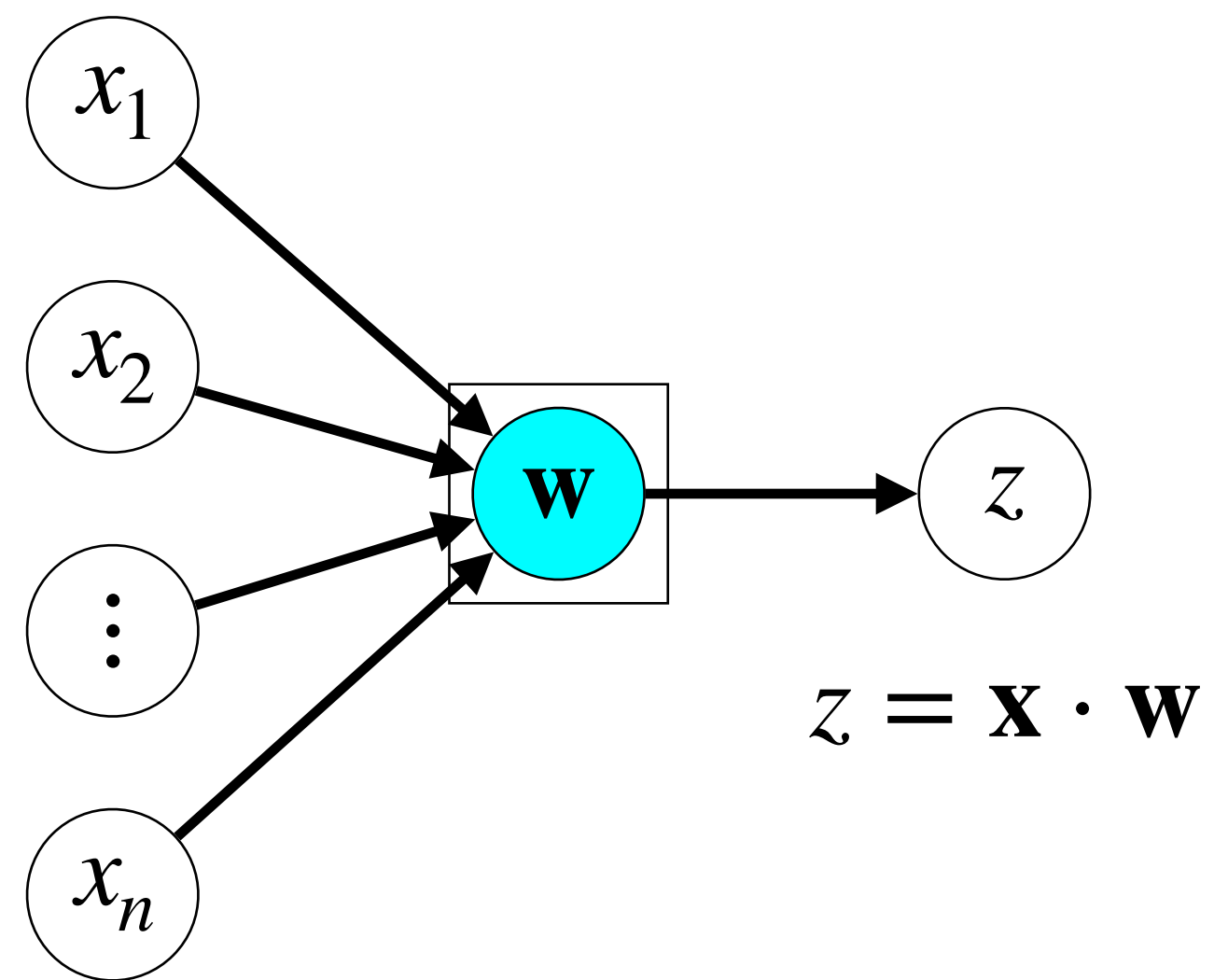
# One neural computation unit (of a hidden layer)



# One neural computation unit (of a hidden layer)

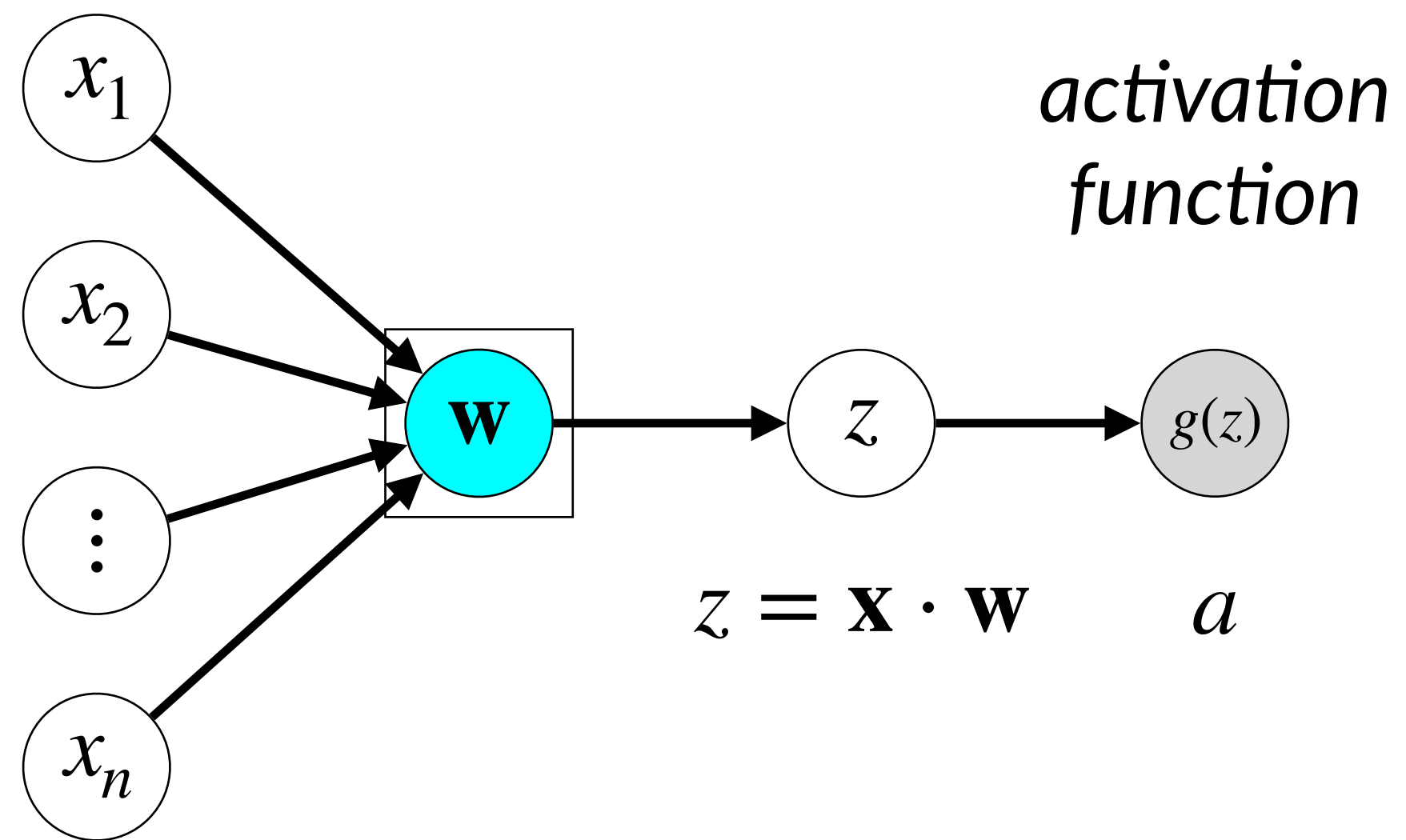


# Activation functions

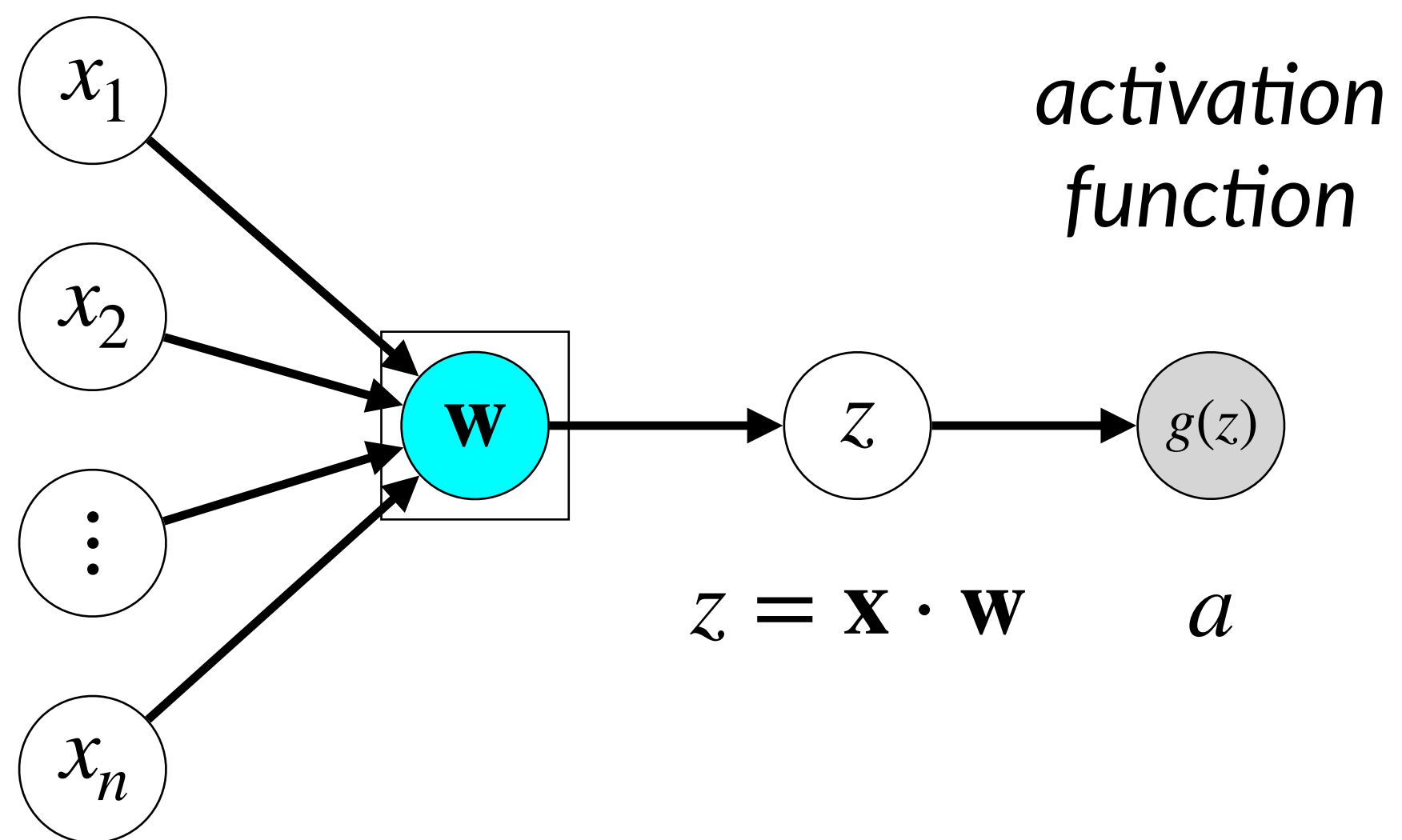




# Activation functions



# Activation functions



$$\sigma(z) = \frac{1}{1 + \exp(-z)} = a$$

*sigmoid  
logistic*

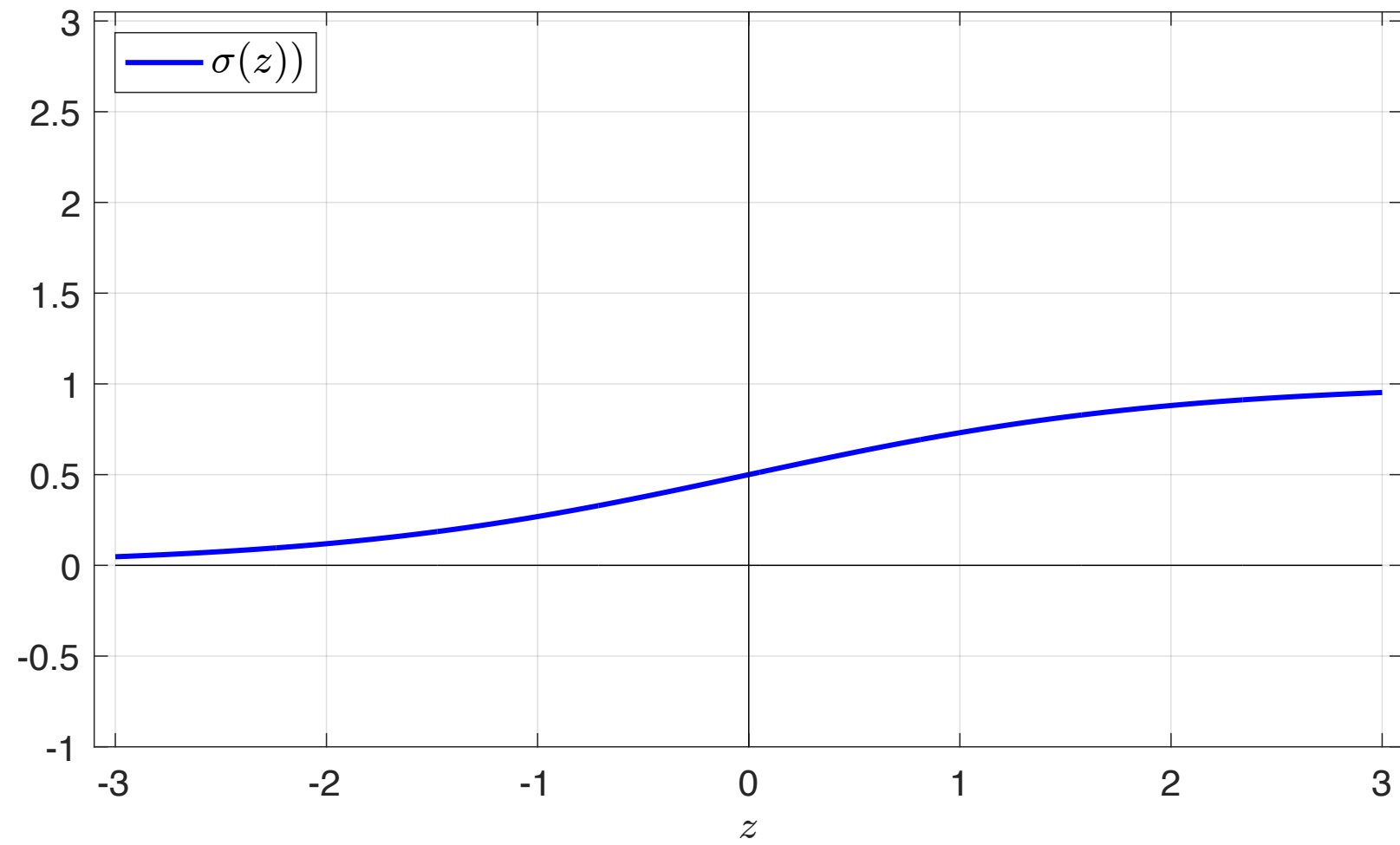
$$\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} = a$$

*hyperbolic  
tangent*

$$\text{ReLU}(z) = \max(z, 0) = a$$

*rectified  
linear unit*

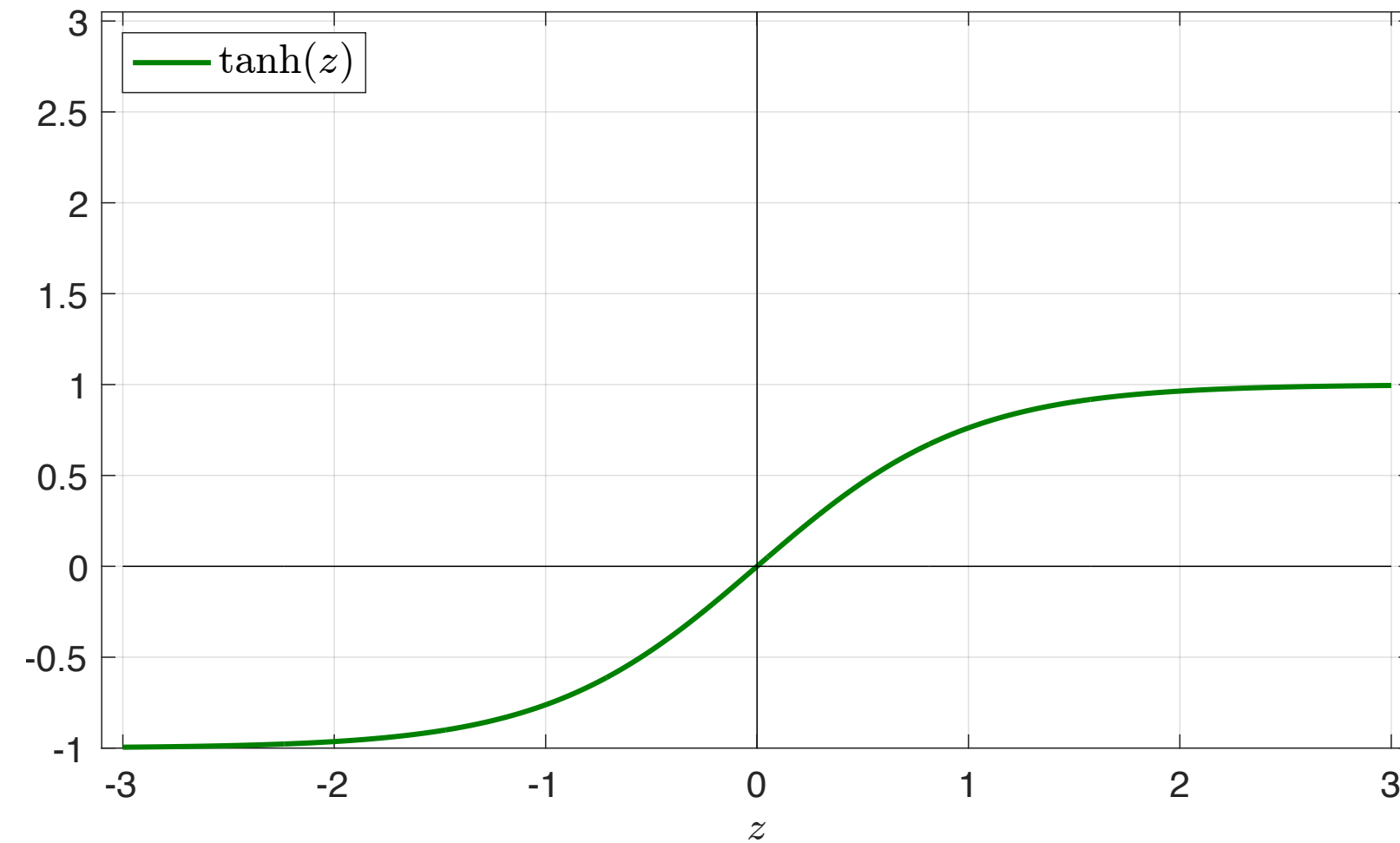
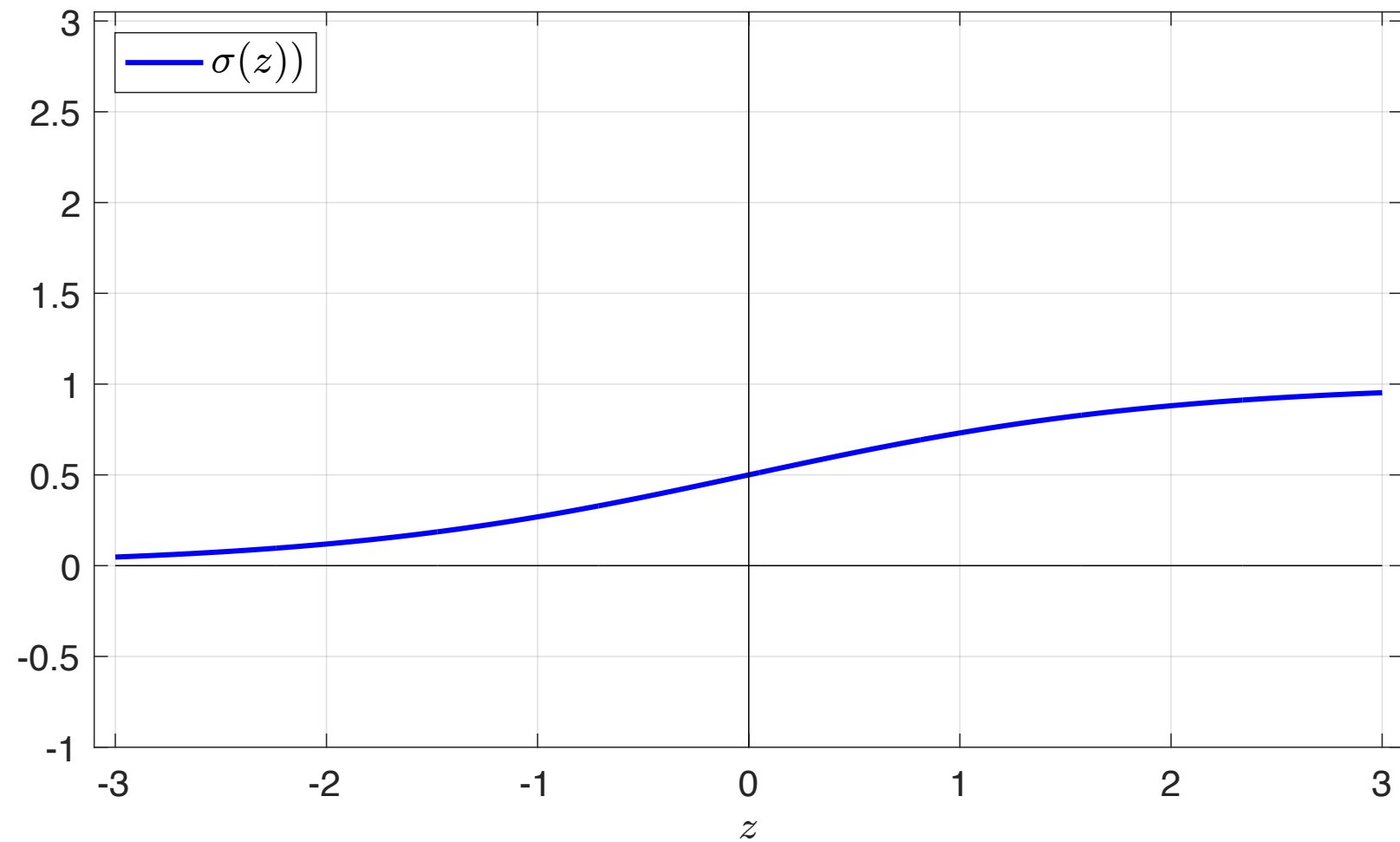
# Activation functions



$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$\sigma \in (0,1)$$

# Activation functions



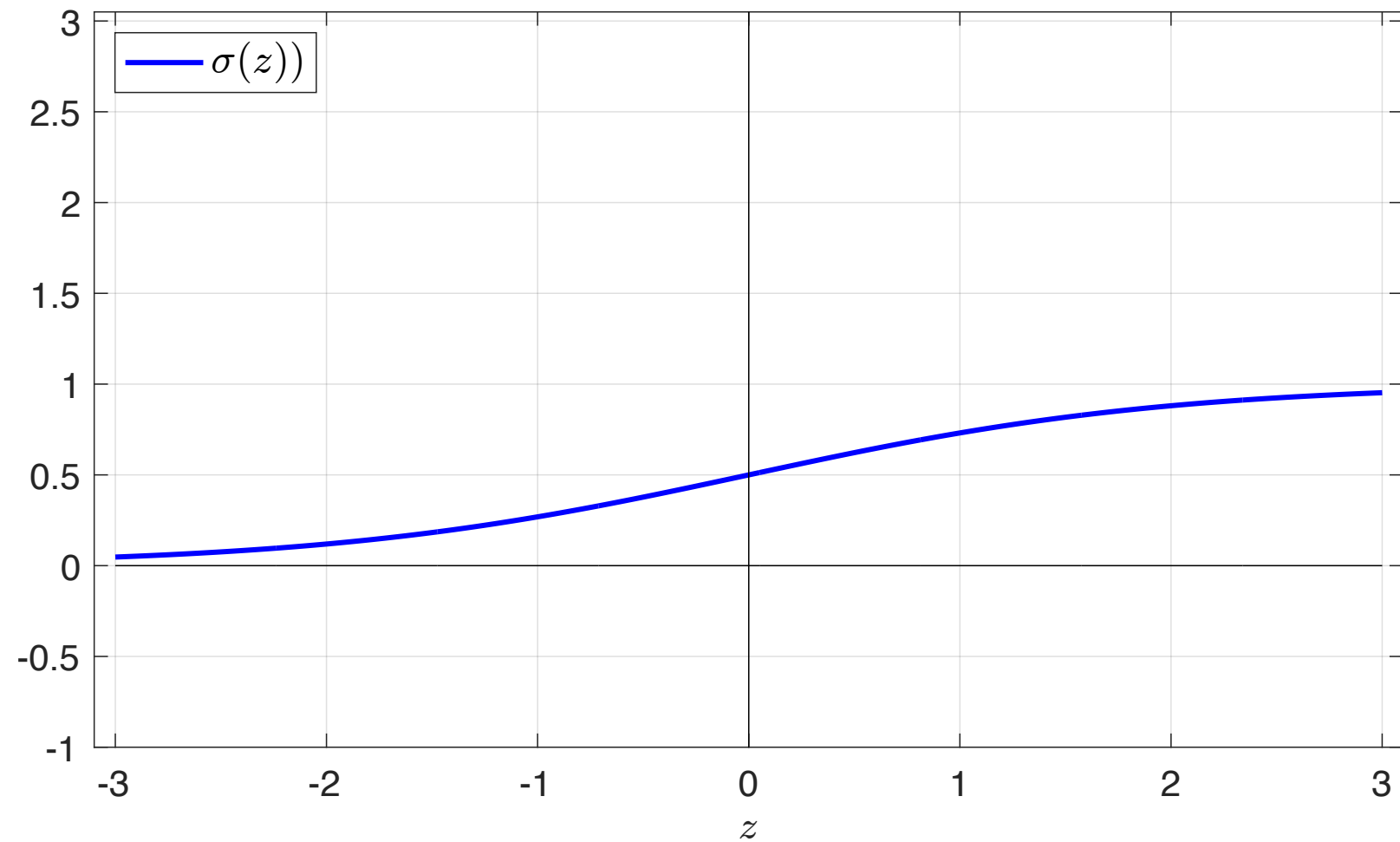
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$\sigma \in (0, 1)$$

$$\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

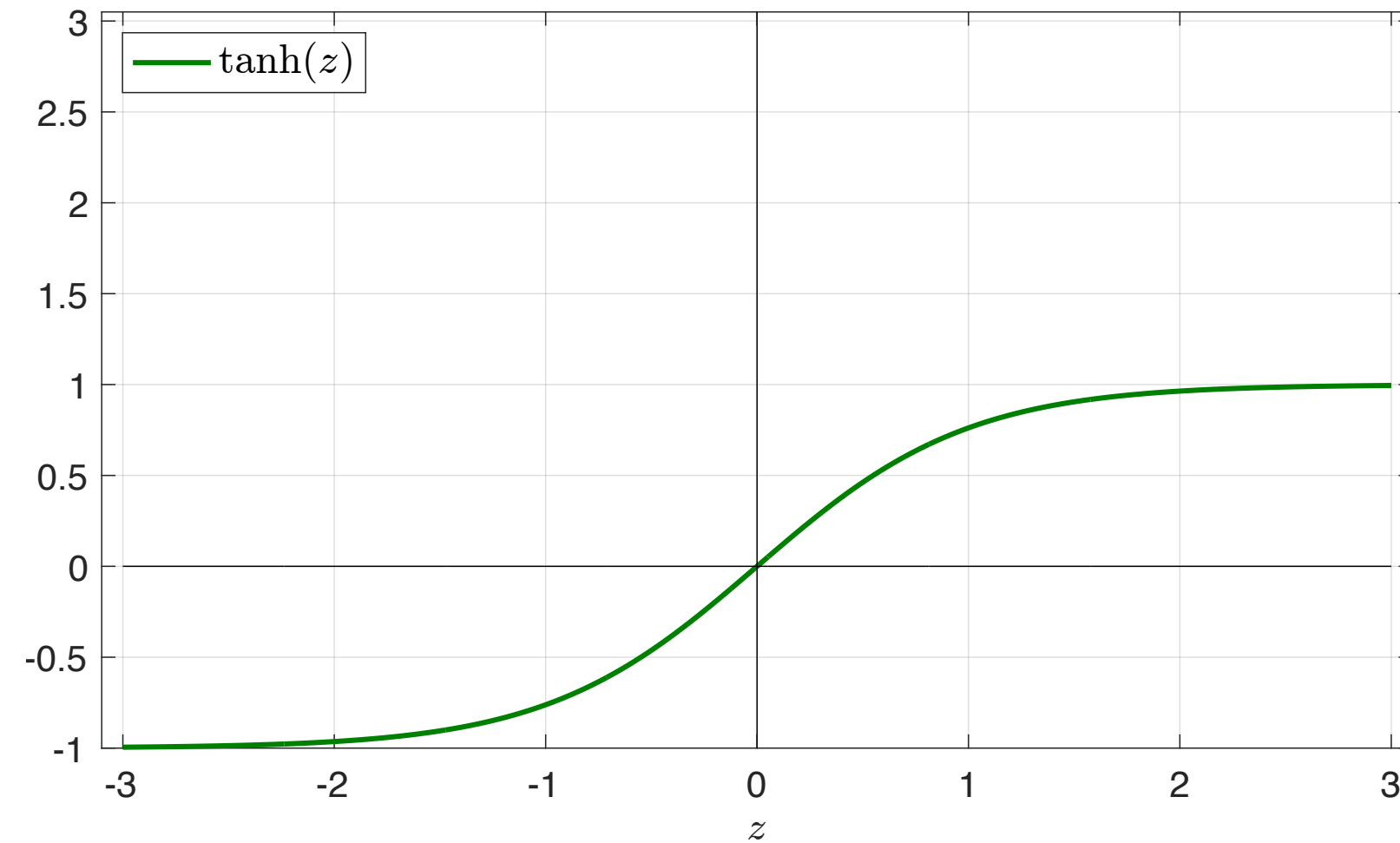
$$\tanh \in (-1, 1)$$

# Activation functions



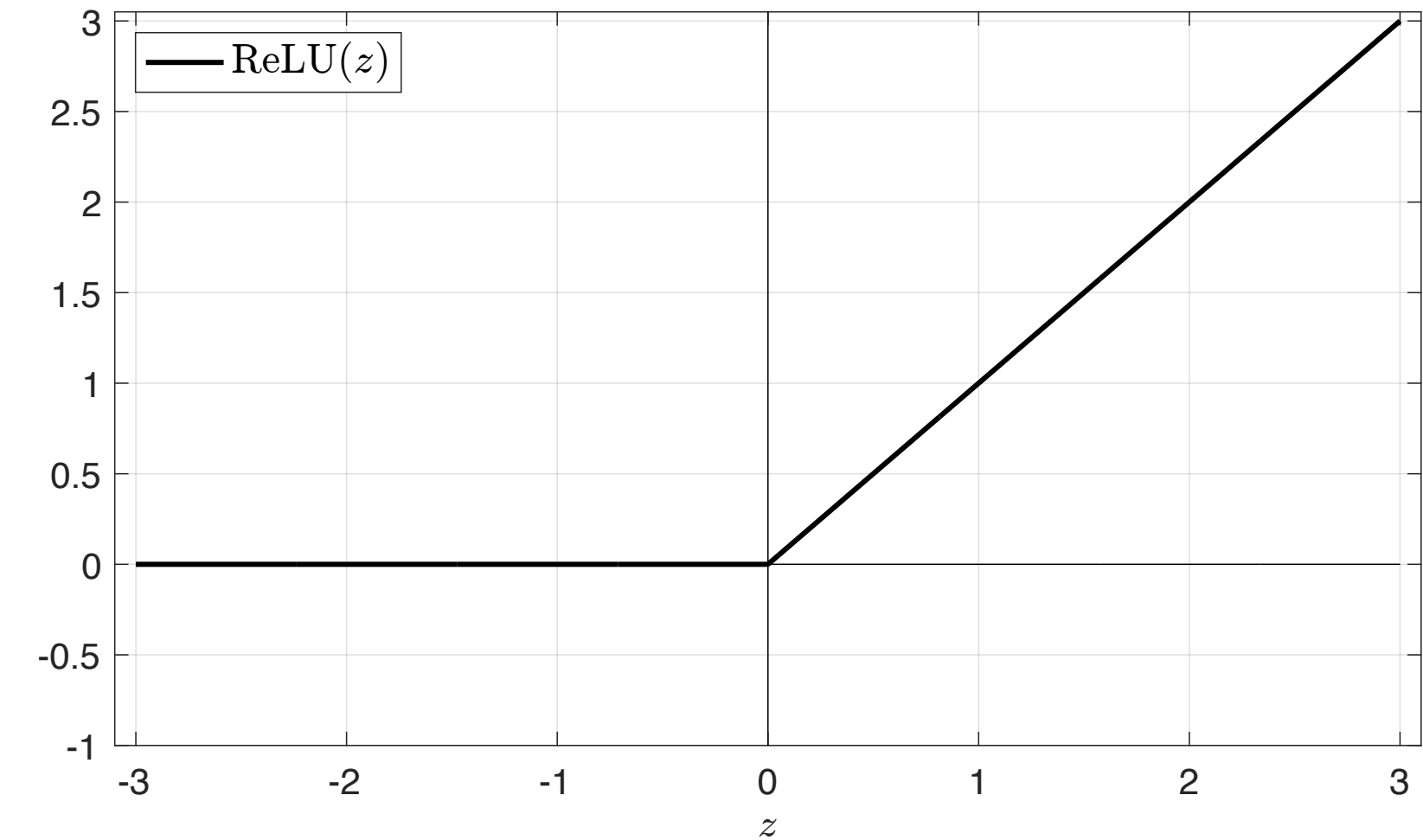
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$\sigma \in (0, 1)$$



$$\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

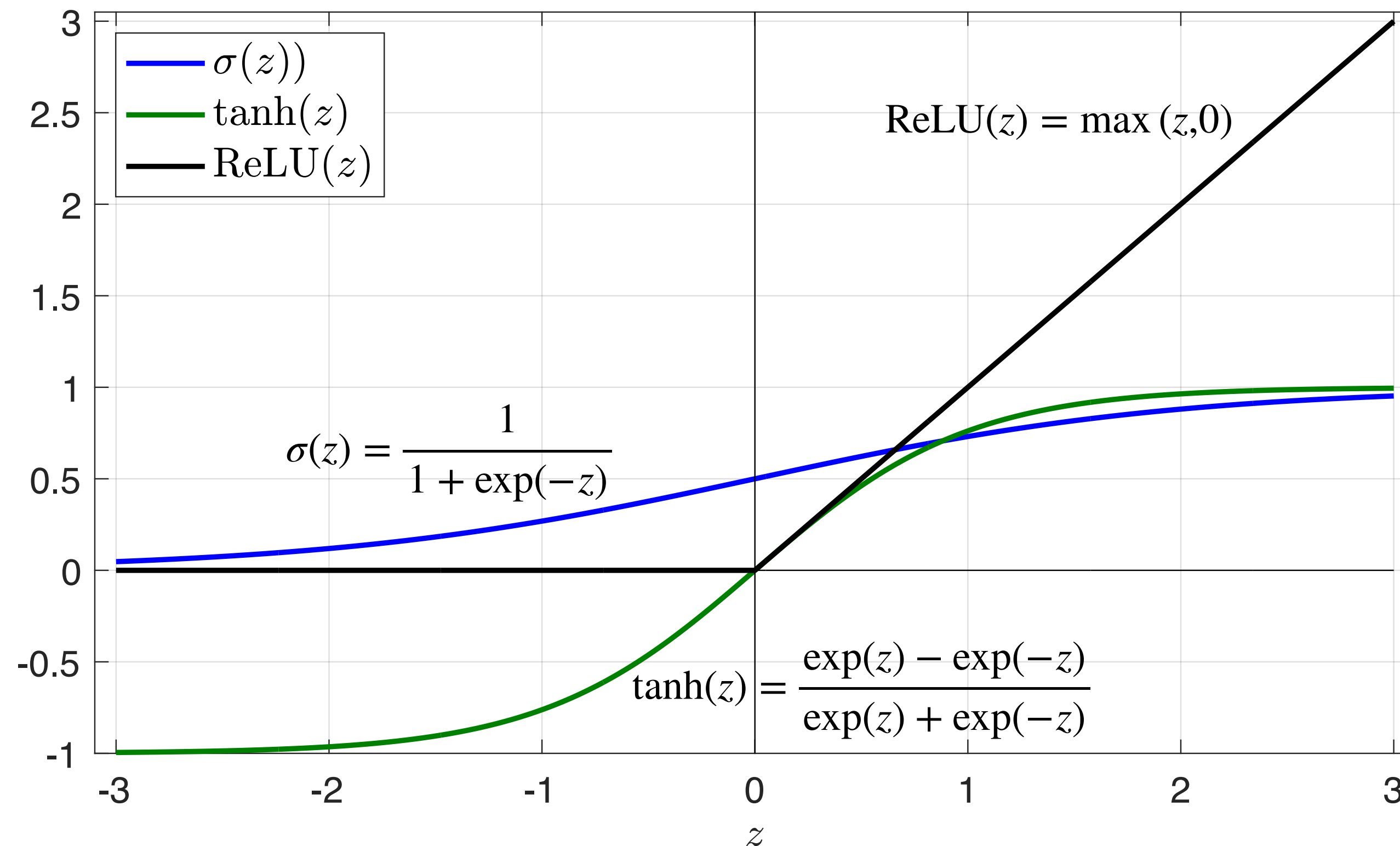
$$\tanh \in (-1, 1)$$



$$\text{ReLU}(z) = \max(z, 0)$$

$$\text{ReLU} \in (0, +\infty)$$

# Activation functions – Vanishing gradient



$\sigma \in (0, 1)$

$\tanh \in (-1, 1)$

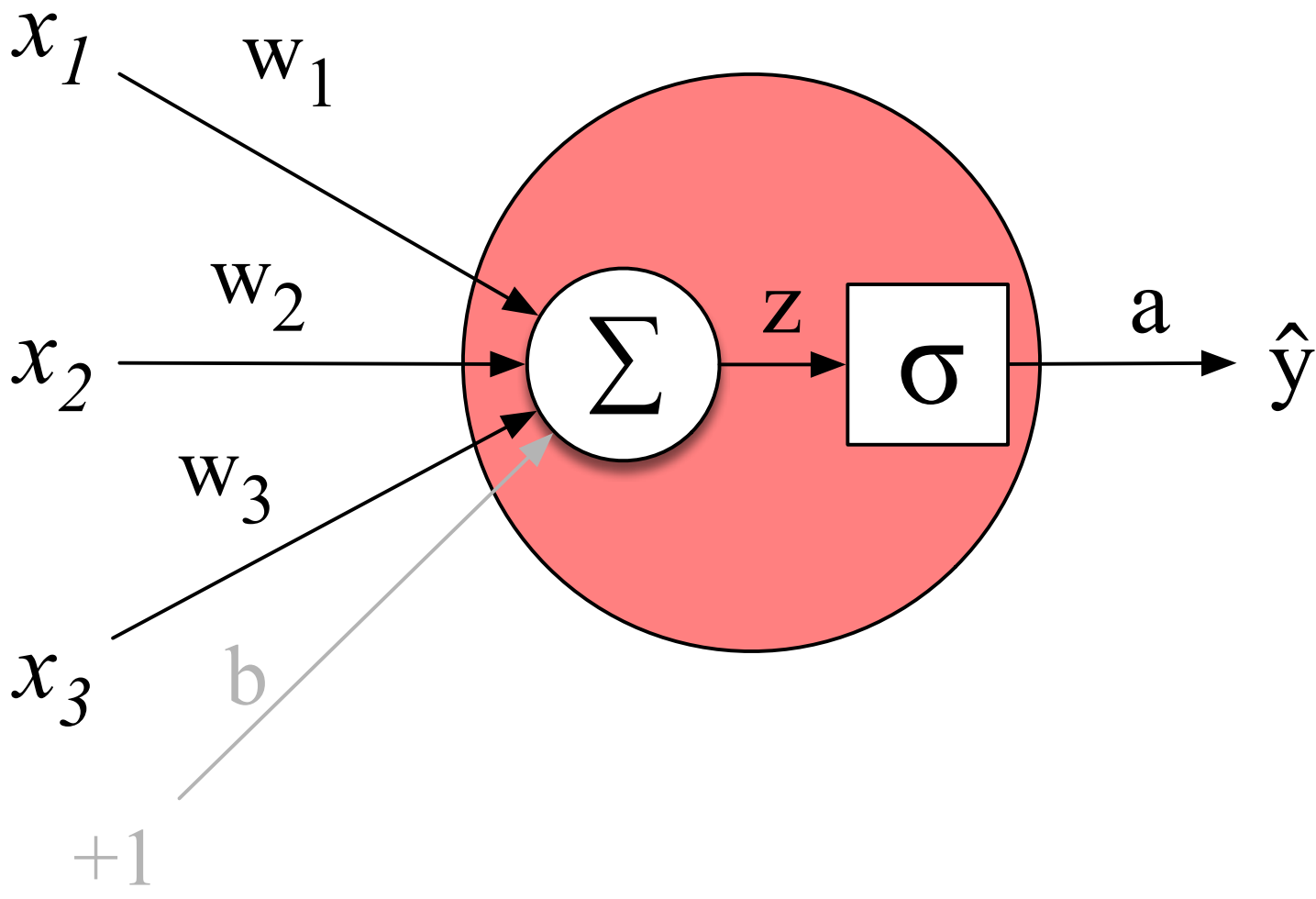
$\text{ReLU} \in (0, +\infty)$

- ▶  $\sigma$ ,  $\tanh$  are differentiable,  $\text{ReLU}$  not differentiable at 0
- ▶  $\tanh$  is almost always preferred to  $\sigma$ , more expansive mapping
- ▶ if  $z \gg 0$ ,  $\sigma$  and  $\tanh$  become saturated, i.e.  $\approx 1$  with derivatives  $\approx 0 \rightarrow$  gradient updates  $\approx 0$  (*no more learning*), vanishing gradient issue
- ▶  $\text{ReLU} \sim$  linear / does not have this vanishing gradient issue

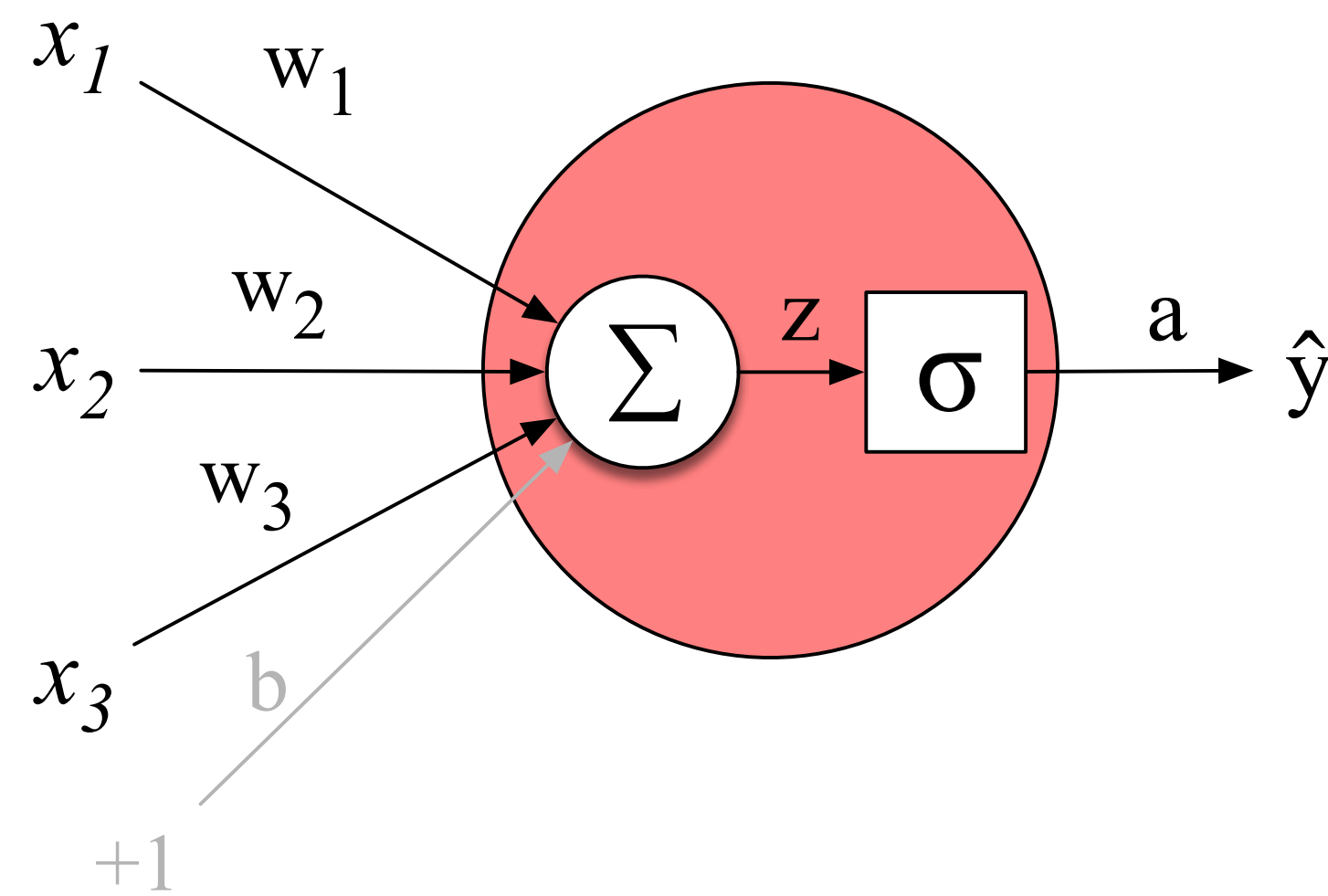


# One neural unit – Example

$$\mathbf{x} = [0.5 \ 0.6 \ 0.1]$$
$$\mathbf{w} = [0.2 \ 0.3 \ 0.9]$$
$$b = 0.5$$



# One neural unit – Example



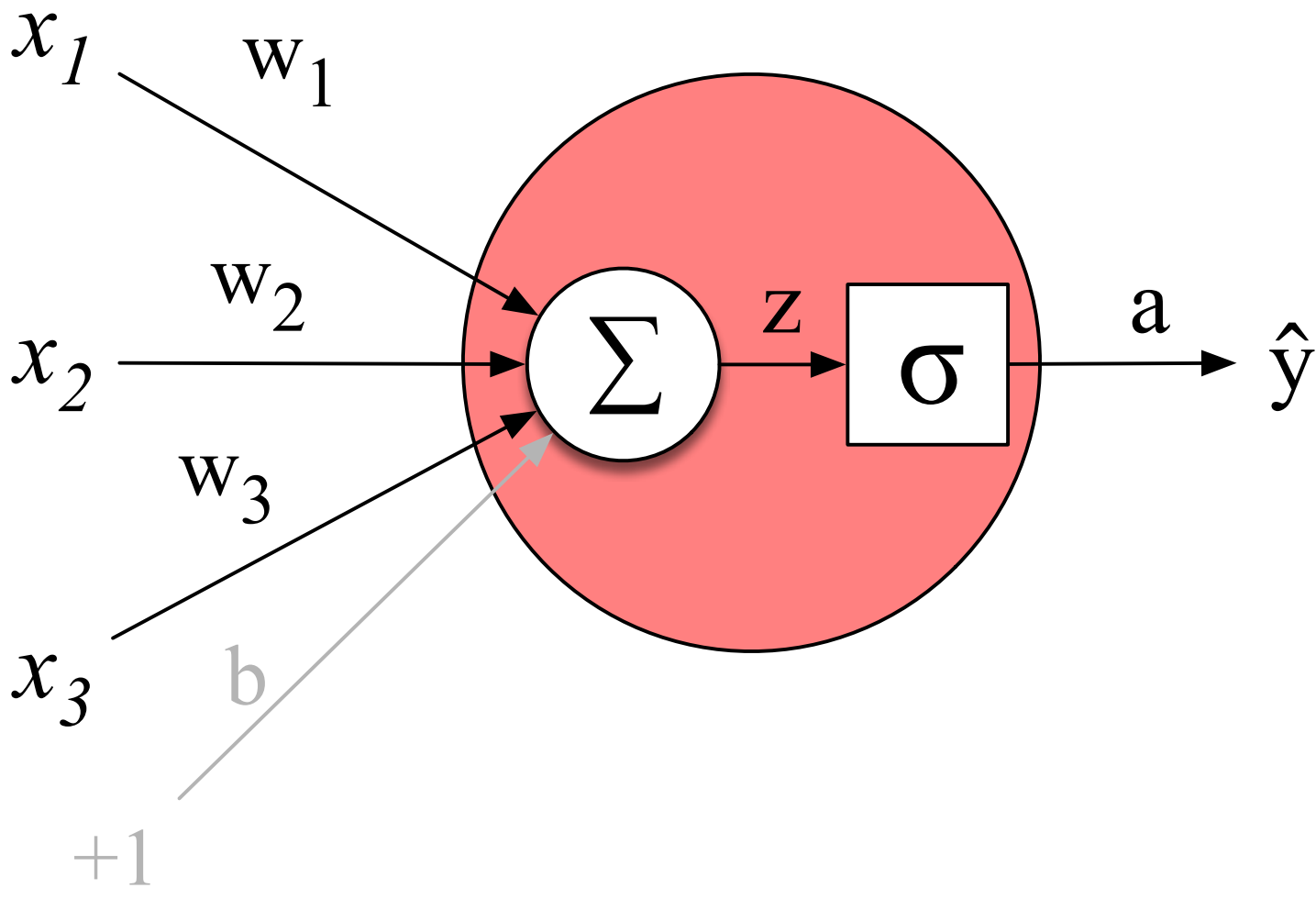
$$\mathbf{x} = [0.5 \ 0.6 \ 0.1]$$

$$\mathbf{w} = [0.2 \ 0.3 \ 0.9]$$

$$b = 0.5$$

$$z = ?$$

# One neural unit – Example



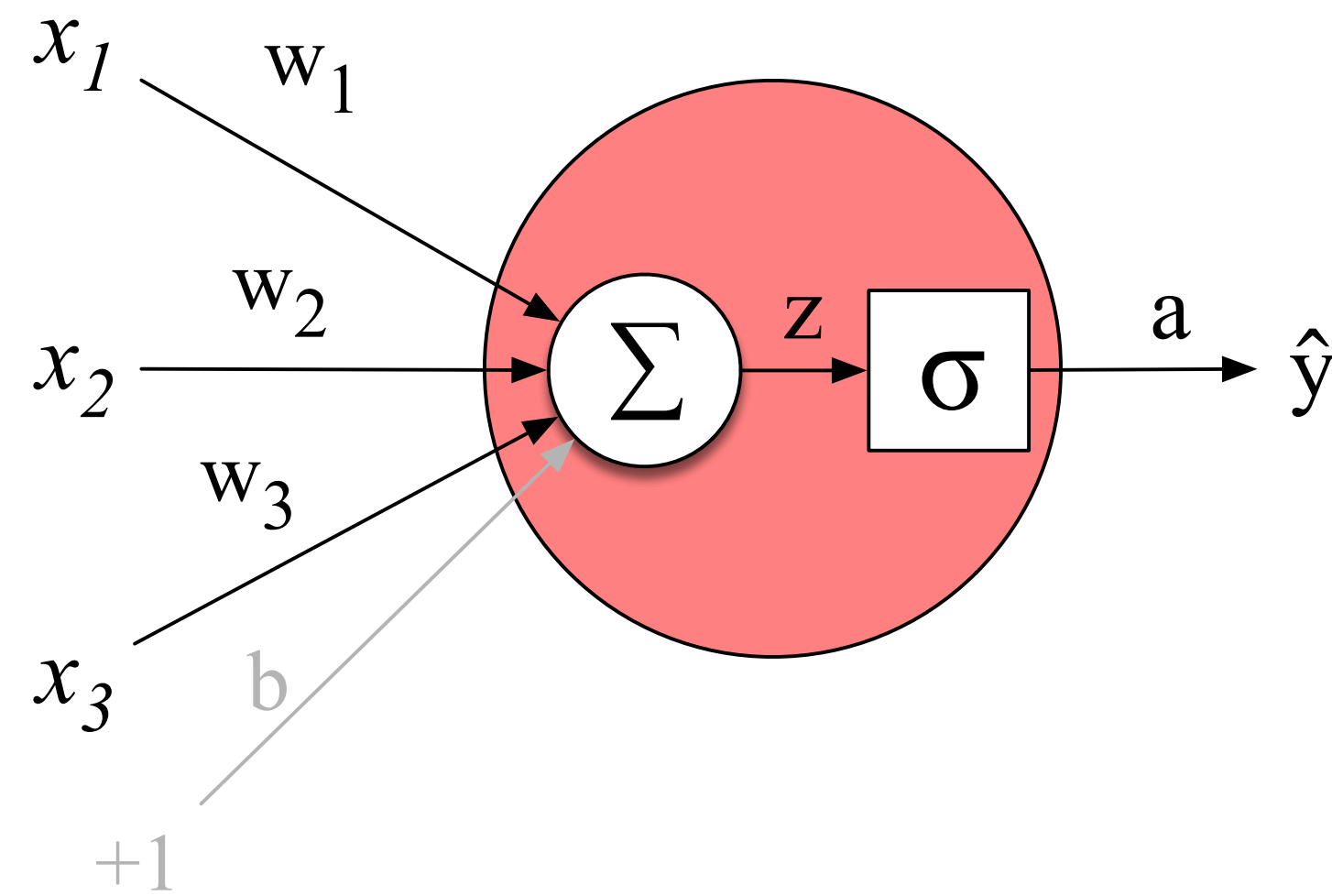
$$\mathbf{x} = [0.5 \ 0.6 \ 0.1]$$
$$\mathbf{w} = [0.2 \ 0.3 \ 0.9]$$
$$b = 0.5$$

----->  
bias

$$\mathbf{x} = [1 \ 0.5 \ 0.6 \ 0.1]$$
$$\mathbf{w} = [0.5 \ 0.2 \ 0.3 \ 0.9]$$

$z = ?$

# One neural unit – Example



$$\mathbf{x} = [0.5 \ 0.6 \ 0.1]$$

$$\mathbf{w} = [0.2 \ 0.3 \ 0.9]$$

$$b = 0.5$$

----->  
bias

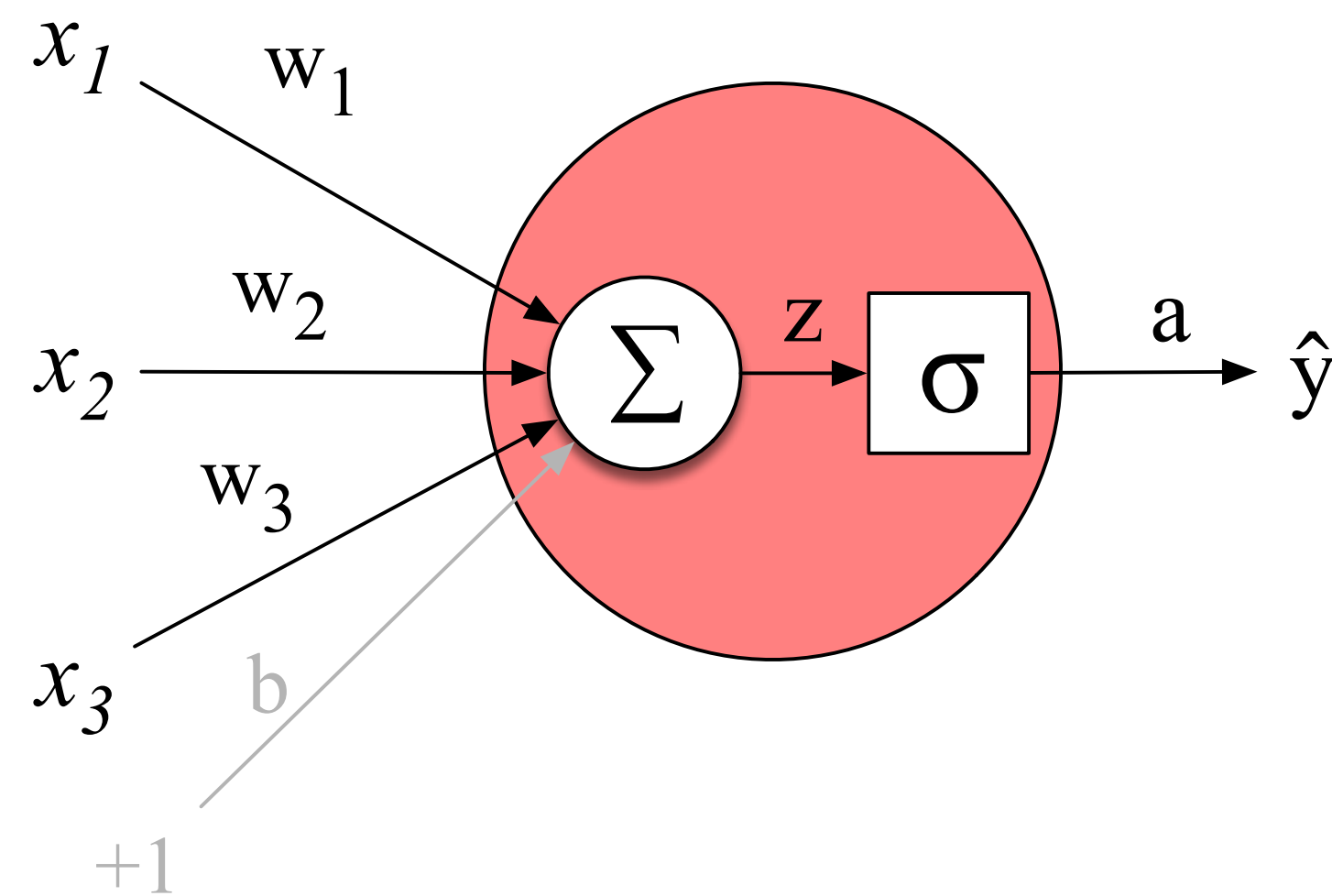
$$\mathbf{x} = [1 \ 0.5 \ 0.6 \ 0.1]$$

$$\mathbf{w} = [0.5 \ 0.2 \ 0.3 \ 0.9]$$

$$z = ?$$

$$z = \mathbf{x} \cdot \mathbf{w} = 1 \cdot 0.5 + 0.5 \cdot 0.2 + \dots + 0.1 \cdot 0.9 = 0.87$$

# One neural unit – Example



$$\mathbf{x} = [0.5 \ 0.6 \ 0.1]$$

$$\mathbf{w} = [0.2 \ 0.3 \ 0.9]$$

$$b = 0.5$$

----->  
bias

$$\mathbf{x} = [1 \ 0.5 \ 0.6 \ 0.1]$$

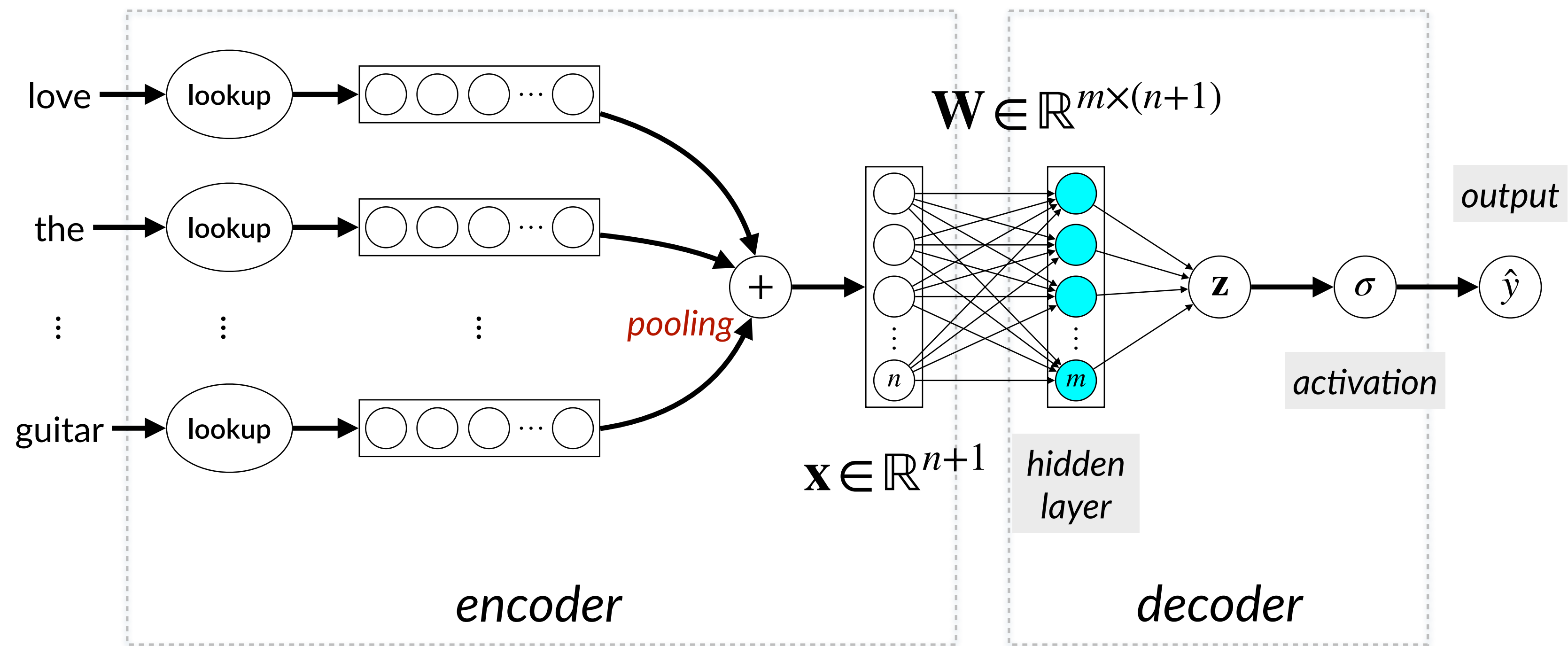
$$\mathbf{w} = [0.5 \ 0.2 \ 0.3 \ 0.9]$$

$$z = ?$$

$$z = \mathbf{x} \cdot \mathbf{w} = 1 \cdot 0.5 + 0.5 \cdot 0.2 + \dots + 0.1 \cdot 0.9 = 0.87$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + \exp(-z)} = \frac{1}{1 + \exp(-0.87)} = 0.705$$

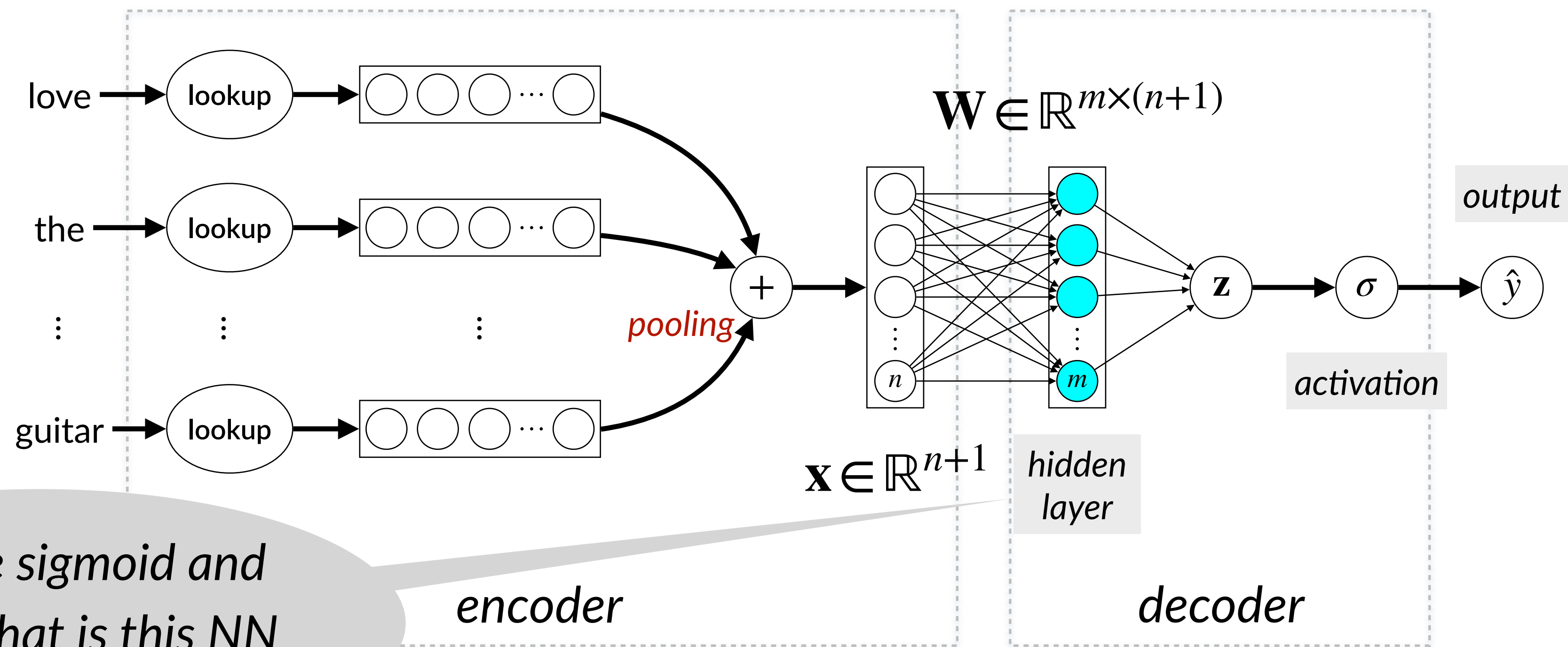
# Feedforward neural network (1 layer)



- A feedforward NN has: input units, **hidden units**, and output units
- Fully connected (*standard version*)
- This NN has 1 layer (*input layer does not count*)



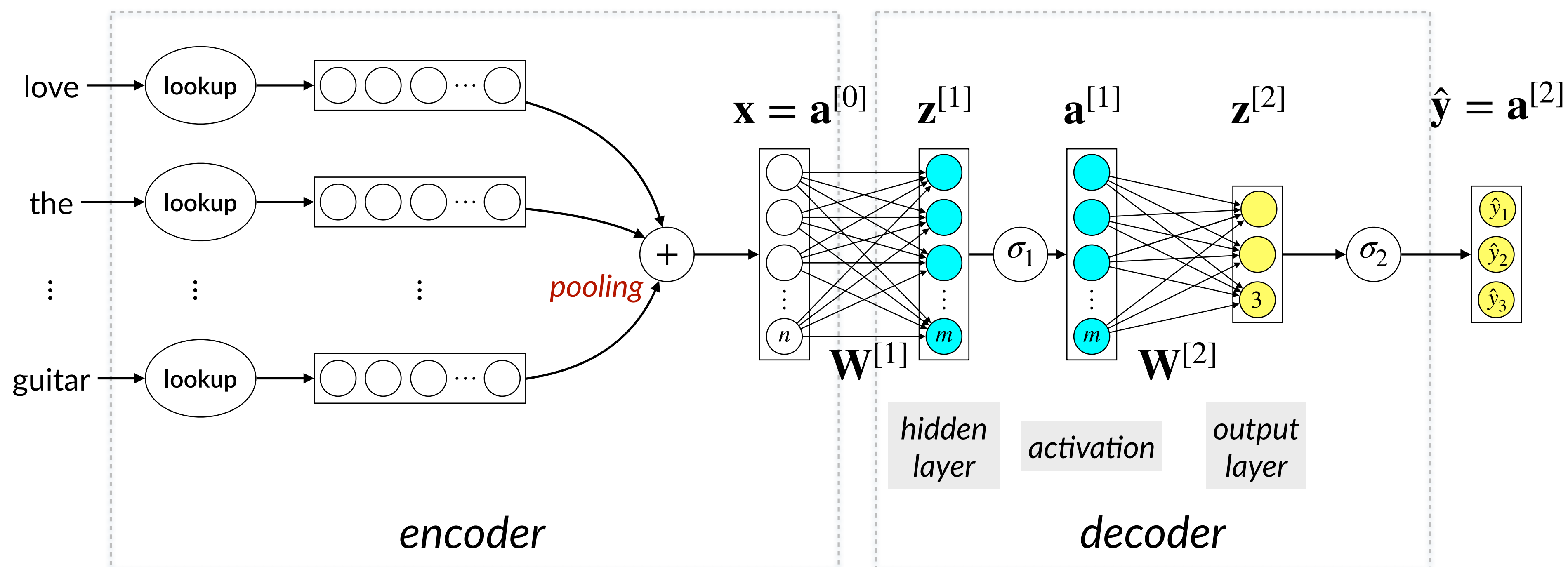
# Feedforward neural network (1 layer)



if  $\sigma$  is the sigmoid and  $m = 1$  what is this NN ~equivalent to?

- A feedforward NN has: input units, **hidden units**, and output units
- Fully connected (*standard version*)
- This NN has 1 layer (*input layer does not count*)

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

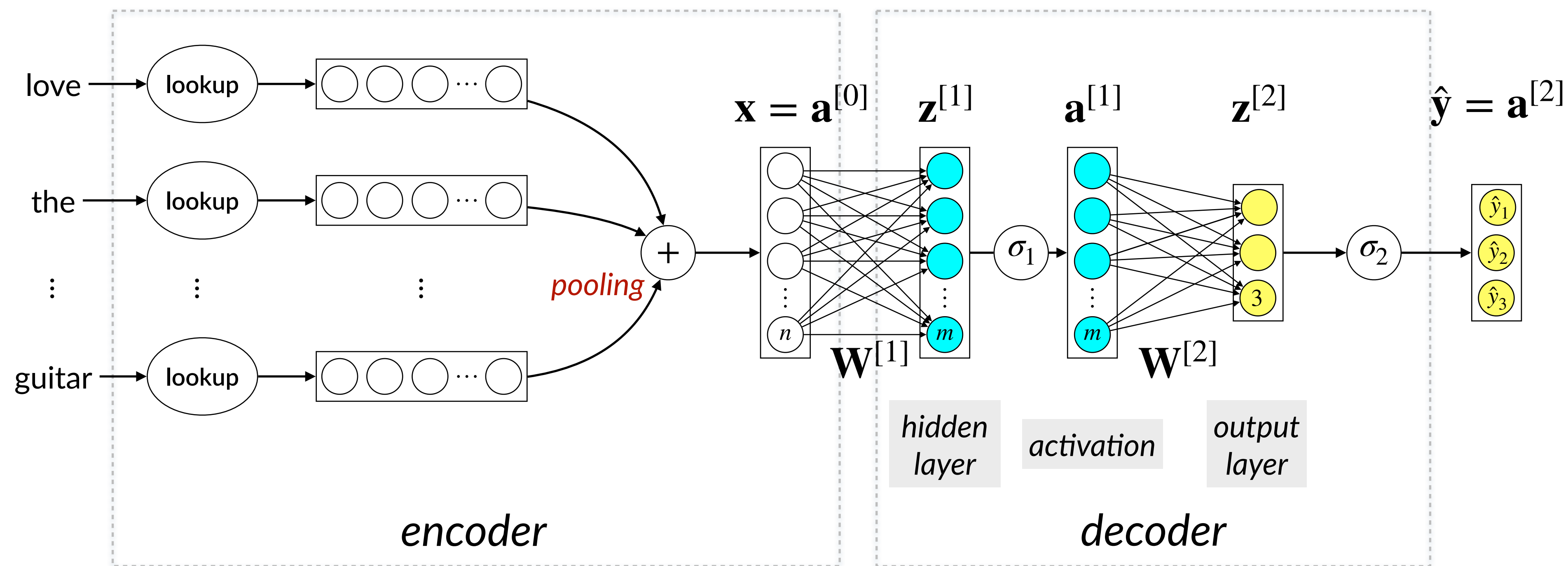
$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

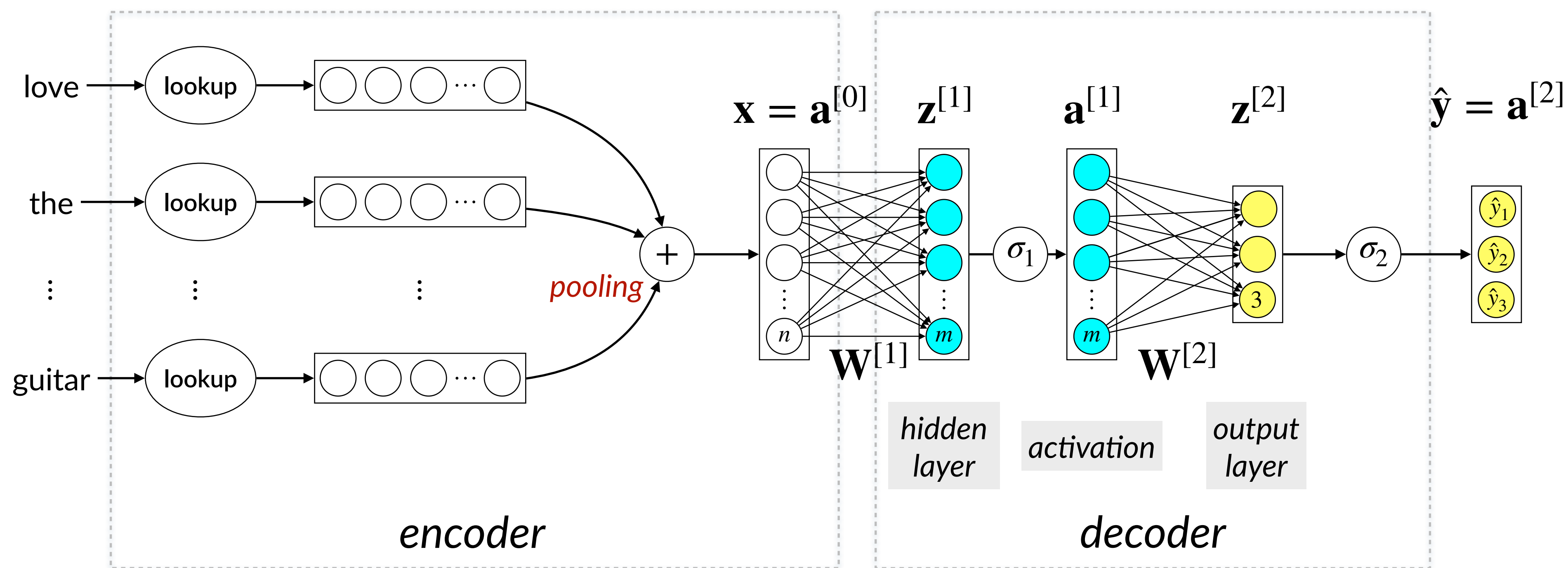
$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

Dimensionalities?

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

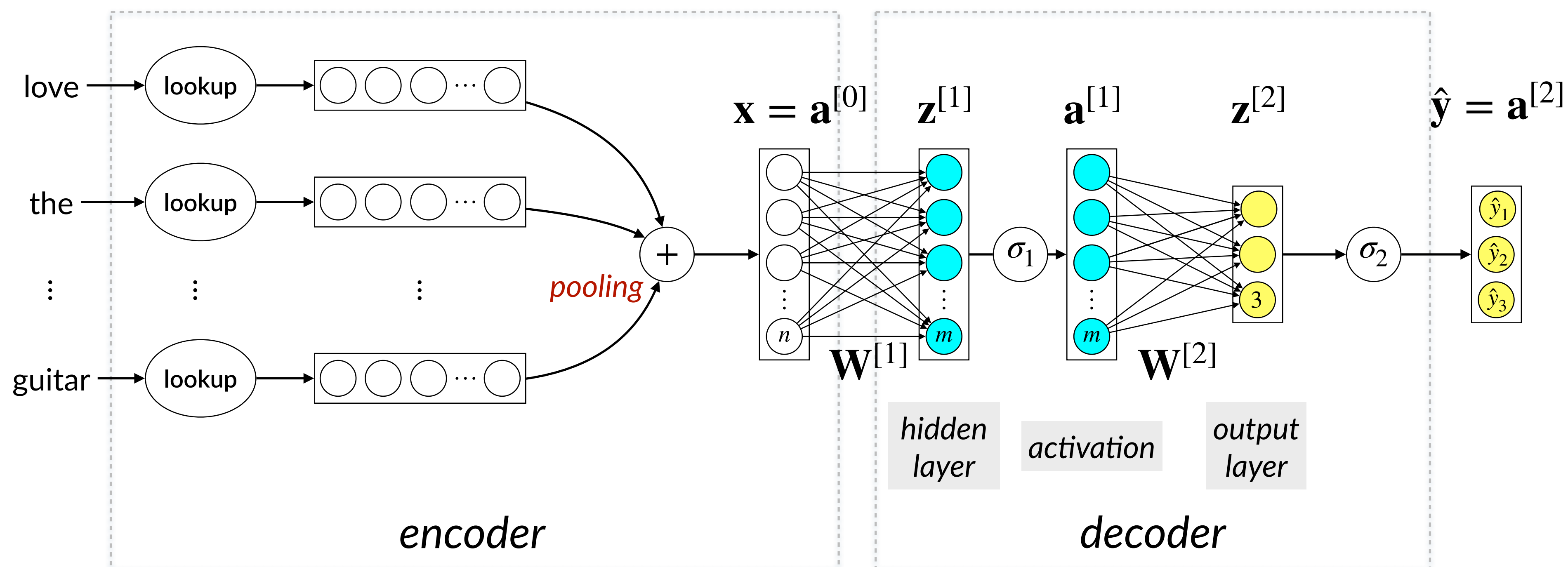
$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

Dimensionalities?



# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

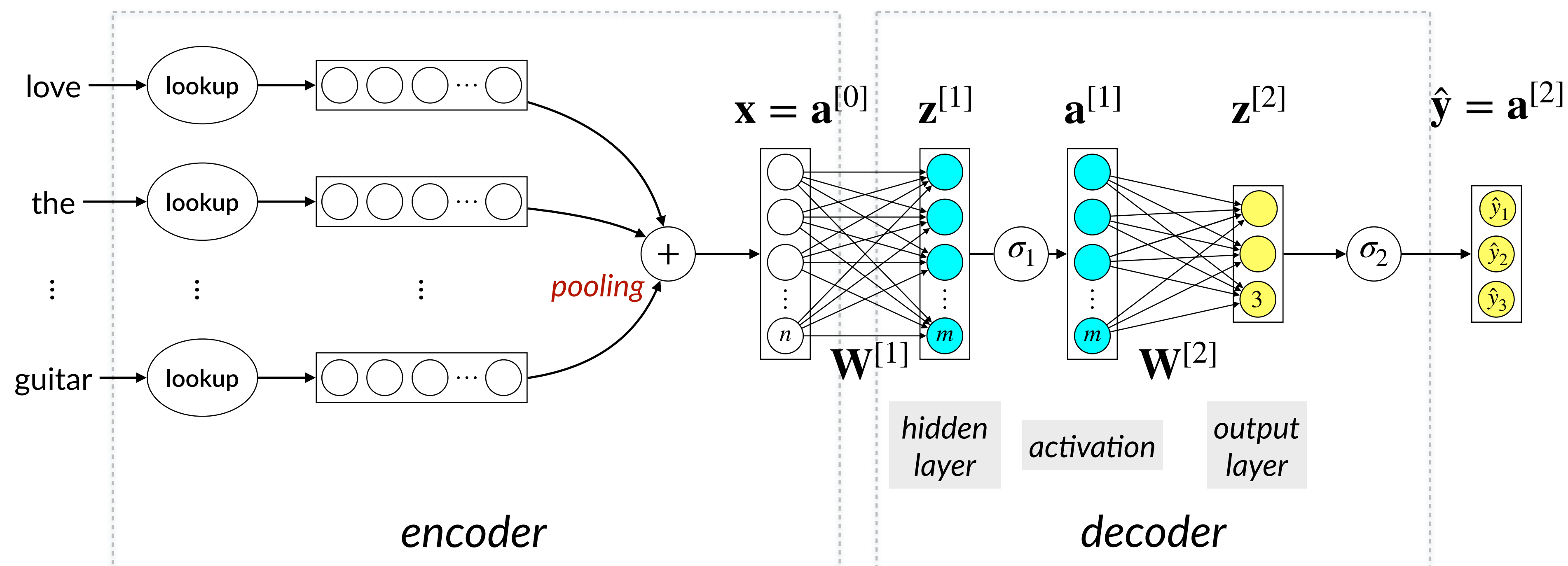
$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

bias term (+1)

Dimensionalities?

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

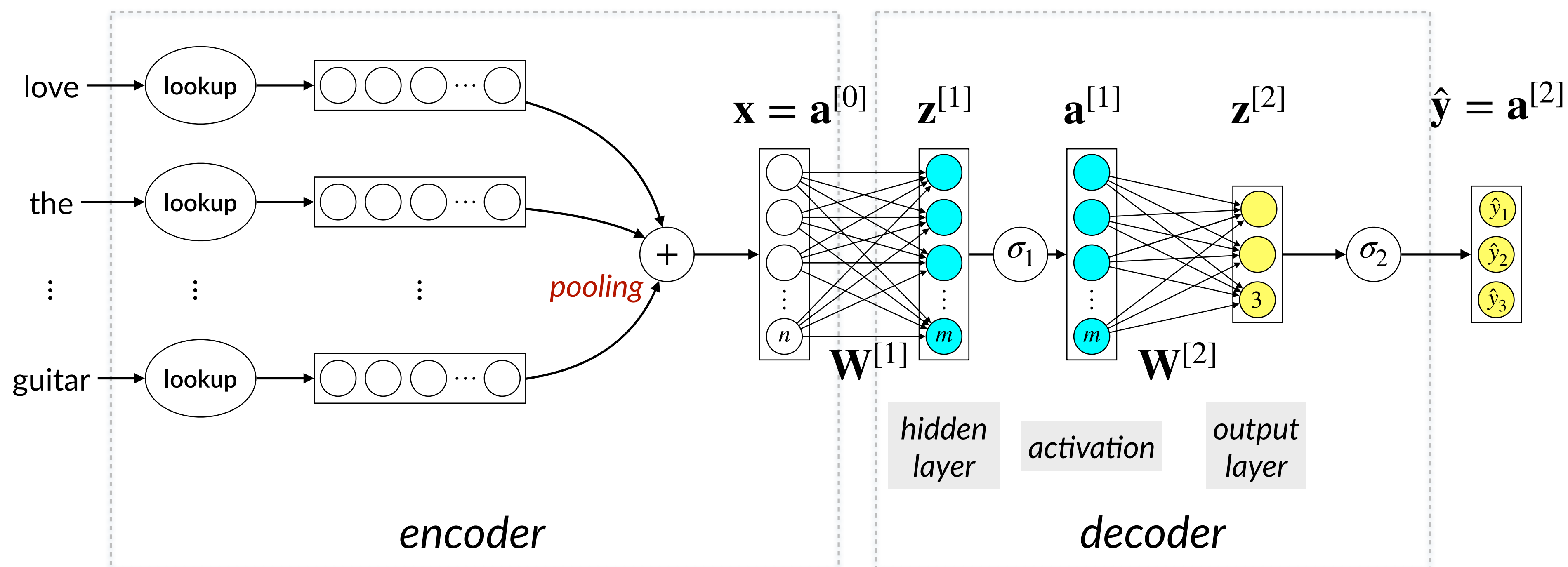
$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

$$\mathbf{a}^{[1]}, \mathbf{z}^{[1]} \in \mathbb{R}^m$$

bias term (+1)

Dimensionalities?

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

$$\mathbf{a}^{[1]}, \mathbf{z}^{[1]} \in \mathbb{R}^m$$

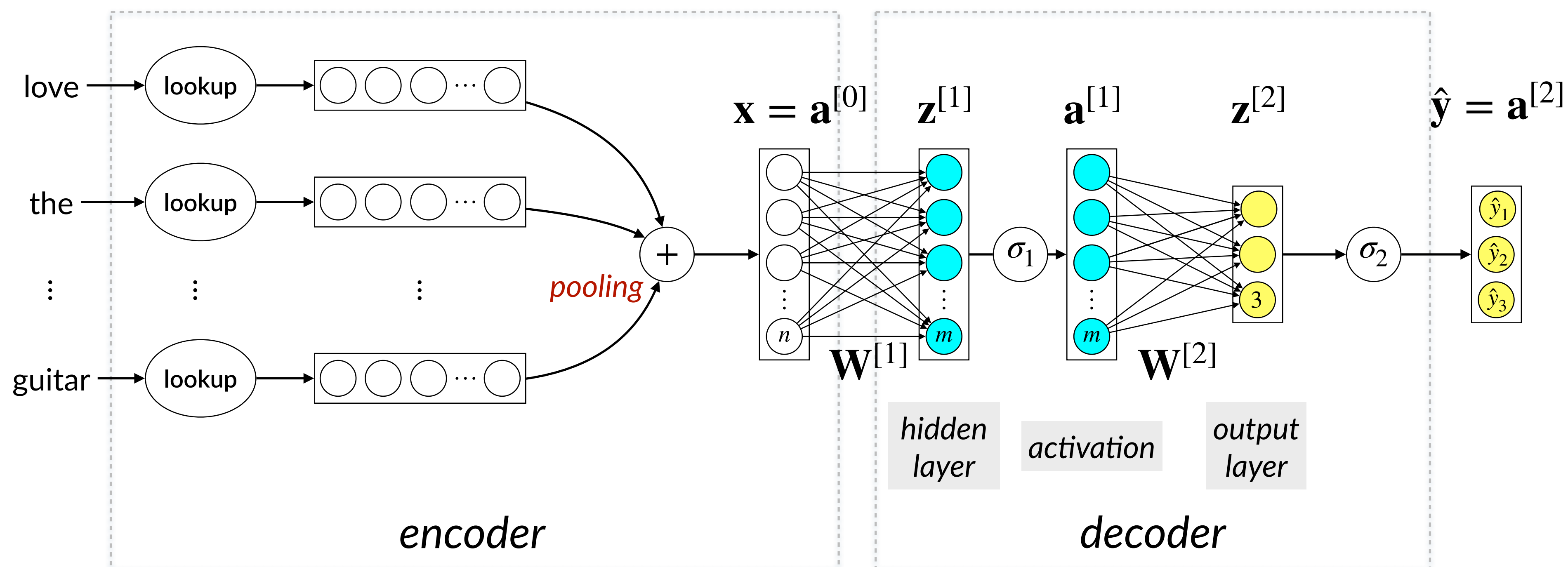
$$\mathbf{W}^{[2]} \in \mathbb{R}^{3 \times m}$$

bias term (+1)

Dimensionalities?



# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

$$\mathbf{a}^{[1]}, \mathbf{z}^{[1]} \in \mathbb{R}^m$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{3 \times m}$$

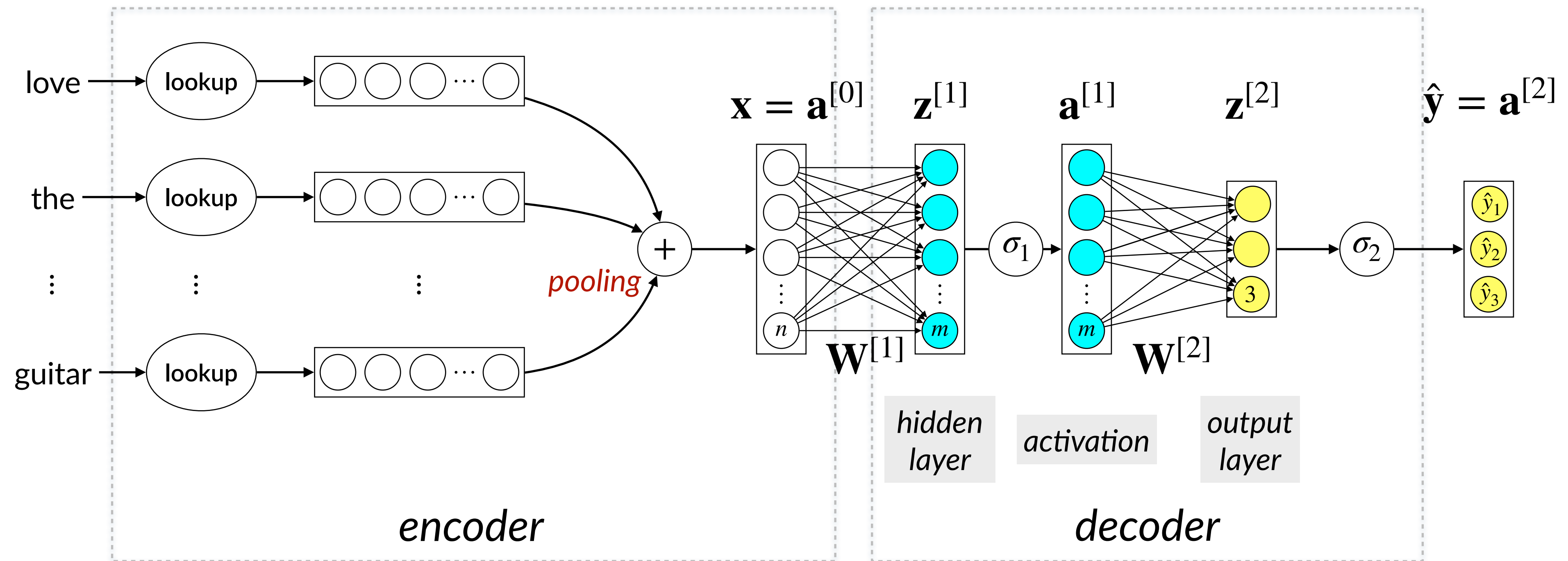
$$\mathbf{a}^{[2]}, \mathbf{z}^{[2]}, \hat{\mathbf{y}} \in \mathbb{R}^3$$

bias term (+1)

no bias term used in the output layer

Dimensionalities?

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

$$\mathbf{a}^{[1]}, \mathbf{z}^{[1]} \in \mathbb{R}^m$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{3 \times m}$$

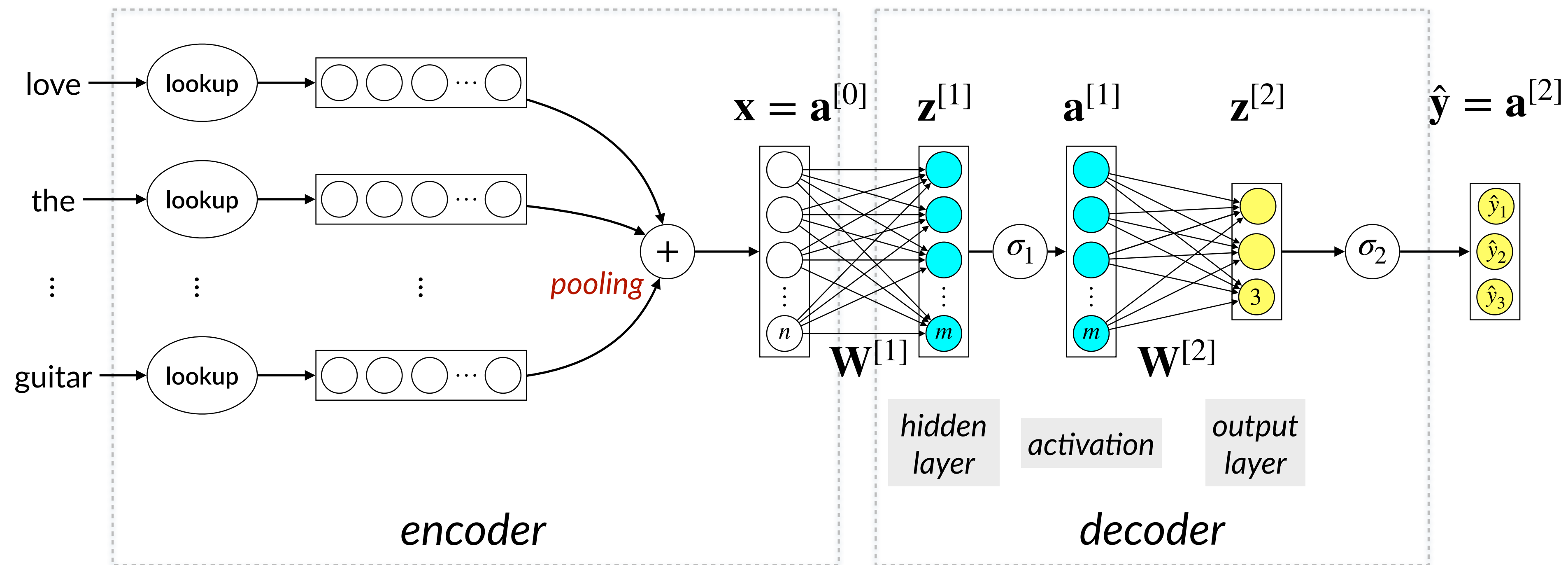
$$\mathbf{a}^{[2]}, \mathbf{z}^{[2]}, \hat{\mathbf{y}} \in \mathbb{R}^3$$

bias term (+1)

no bias term used in  
the output layer

$$\hat{\mathbf{y}} = \sigma_2\left(\mathbf{W}^{[2]} \sigma_1\left(\mathbf{W}^{[1]} \mathbf{x}\right)\right)$$

# Feedforward neural network (2 layers, multiple outputs)



$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$$

$$\mathbf{a}^{[1]} = \sigma_1(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = \sigma_2(\mathbf{z}^{[2]})$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[2]}$$

$$\mathbf{x}, \mathbf{a}^{[0]} \in \mathbb{R}^{n+1}$$

$$\mathbf{W}^{[1]} \in \mathbb{R}^{m \times (n+1)}$$

$$\mathbf{a}^{[1]}, \mathbf{z}^{[1]} \in \mathbb{R}^m$$

$$\mathbf{W}^{[2]} \in \mathbb{R}^{3 \times m}$$

$$\mathbf{a}^{[2]}, \mathbf{z}^{[2]}, \hat{\mathbf{y}} \in \mathbb{R}^3$$

bias term (+1)

no bias term used in the output layer

How many parameters do we need to learn?

# Are nonlinear ( $\sigma$ ) activation functions necessary?

If our activation functions were linear in  $\hat{y} = \sigma_2\left(\mathbf{W}^{[2]} \sigma_1\left(\mathbf{W}^{[1]}\mathbf{x}\right)\right) \dots$

# Are nonlinear ( $\sigma$ ) activation functions necessary?

If our activation functions were linear in  $\hat{y} = \sigma_2\left(\mathbf{W}^{[2]} \sigma_1(\mathbf{W}^{[1]}\mathbf{x})\right) \dots$

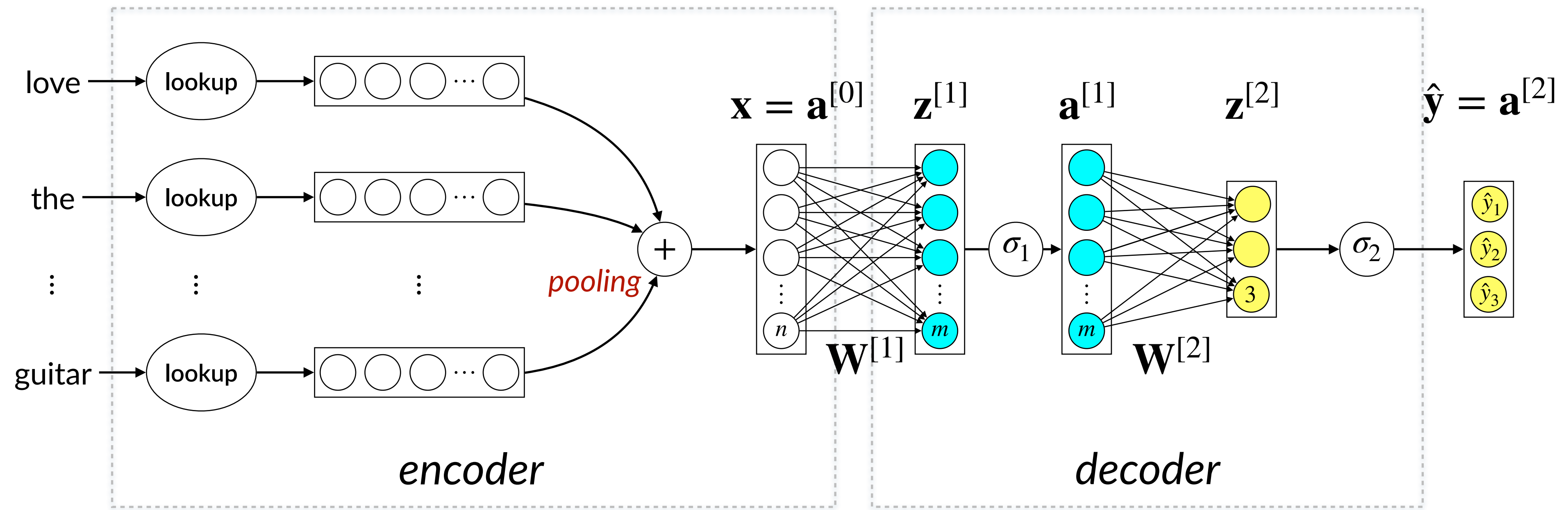
then we can simply omit the non-linear activations  $\sigma_1$  and  $\sigma_2$ :

$$\begin{aligned}\hat{y} &= \mathbf{z}^{[2]} \\ &= \mathbf{W}^{[2]}\mathbf{z}^{[1]} \\ &= \mathbf{W}^{[2]}\mathbf{W}^{[1]}\mathbf{x} \\ &= \mathbf{W}'\mathbf{x}\end{aligned}$$

Hence, we have reduced 2 layers back to 1 with altered parameters ( $\mathbf{W}'$ ).  
This generalises to any number of layers.



# Inference with a feedforward neural network – Softmax



Need to convert outputs to pseudo-probabilities  
➔ common setting for  $\sigma_2$  is the softmax function

$$y_i = \text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^d \exp(z_j)}, \quad 1 \leq i \leq d$$

# Softmax example

$$y_i = \text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^d \exp(z_j)}, \quad 1 \leq i \leq d$$

$$\sum_i y_i = 1 \quad \text{and} \quad y_i \in [0,1] \quad \textit{pseudo-probabilities}$$

So, in our example if

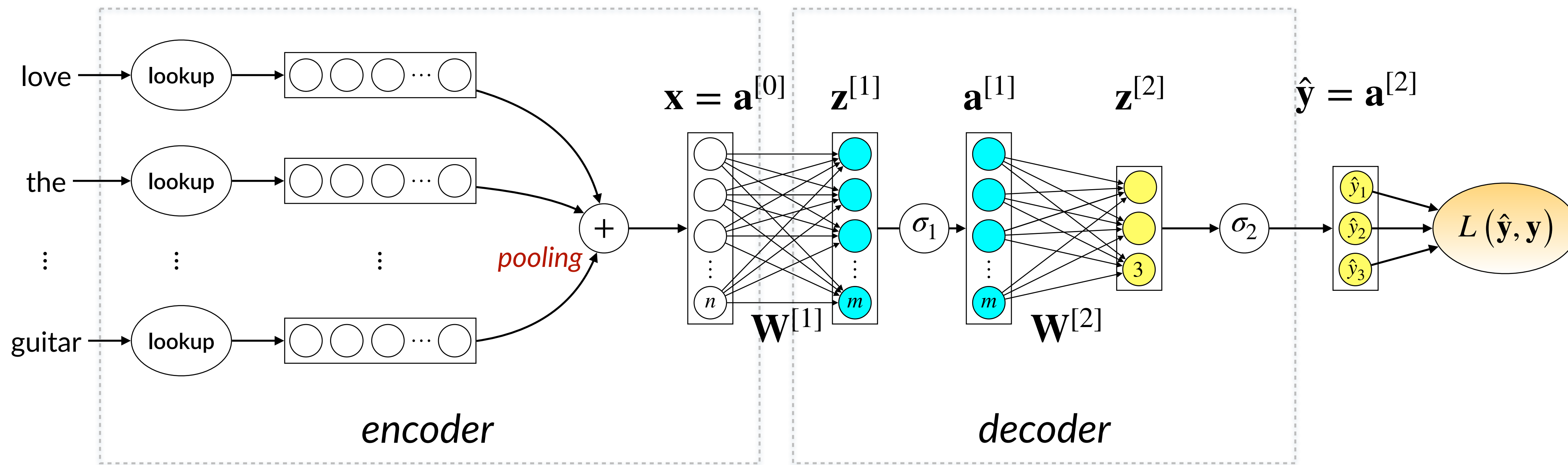
$$\mathbf{z}^{[2]} = [2 \quad -1.99 \quad -0.01]$$

then

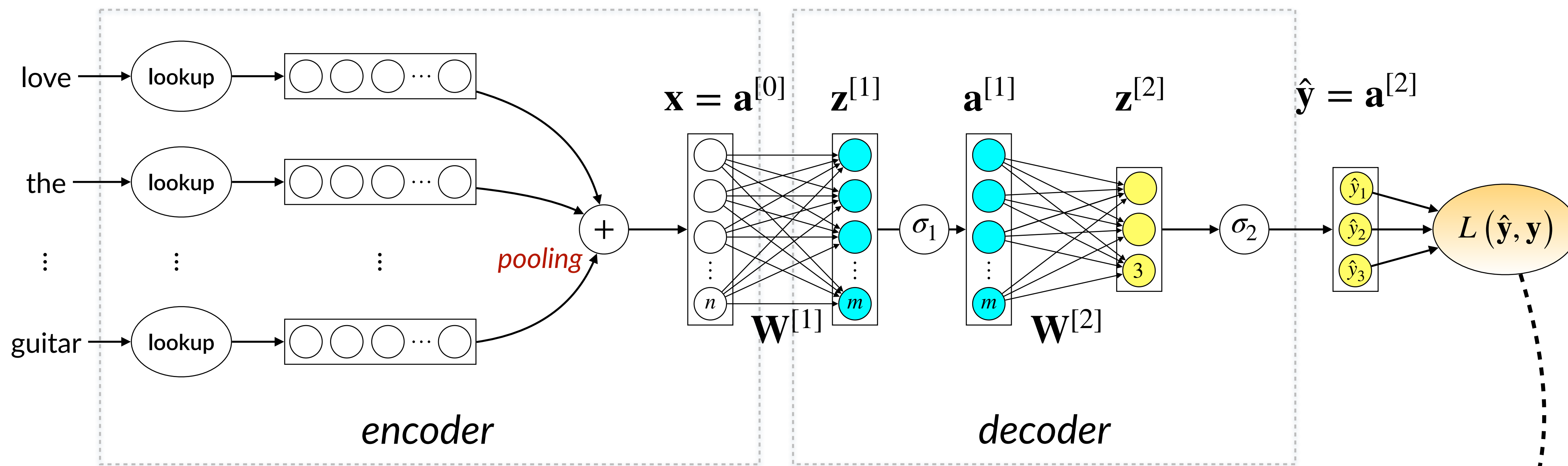
$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z}^{[2]}) = [0.868 \quad 0.016 \quad 0.116]$$



# Training a feedforward neural network

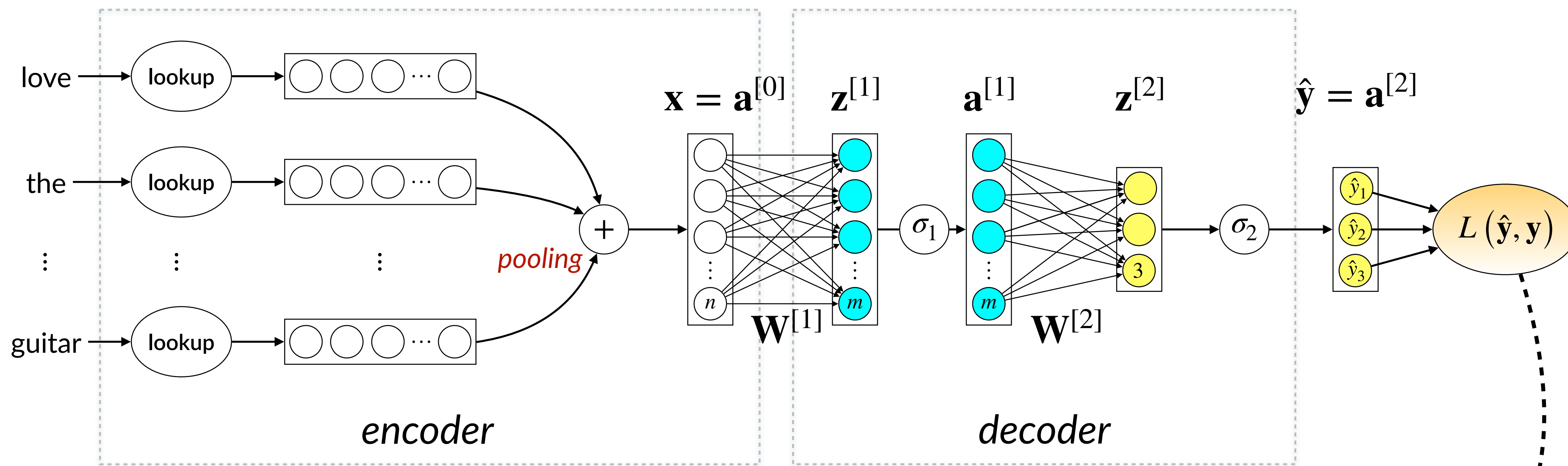


# Training a feedforward neural network



How would I use a loss function  $L(\hat{\mathbf{y}}, \mathbf{y})$  to update *efficiently* my NN parameters in  $\mathbf{W}^{[1]}$  and  $\mathbf{W}^{[2]}$ ?

# Training a feedforward neural network



Backpropagation  
a.k.a.  
"backprop"

How would I use a loss function  $L(\hat{\mathbf{y}}, \mathbf{y})$  to update *efficiently* my NN parameters in  $\mathbf{W}^{[1]}$  and  $\mathbf{W}^{[2]}$ ?

# Cross-entropy loss function

Cross-entropy loss

$$L_{\text{ce}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

where  $K$  is the number of output classes

# Cross-entropy loss function

Cross-entropy loss

$$L_{\text{ce}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

where  $K$  is the number of output classes

Only one of the  $K$   $y_k$ 's will be equal to 1. The rest will be 0.

If, say,  $y_c = 1$ ,  $c = \{1, \dots, K\}$ , i.e.  $c$  is the correct class, the loss can be simplified as:

$$L_{\text{ce}}(\hat{\mathbf{y}}, \mathbf{y}) = - y_c \cdot \log \hat{y}_c = - \log \hat{y}_c$$

# Cross-entropy loss function

Cross-entropy loss

$$L_{\text{ce}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

where  $K$  is the number of output classes

Only one of the  $K$   $y_k$ 's will be equal to 1. The rest will be 0.

If, say,  $y_c = 1$ ,  $c = \{1, \dots, K\}$ , i.e.  $c$  is the correct class, the loss can be simplified as:

$$L_{\text{ce}}(\hat{\mathbf{y}}, \mathbf{y}) = - y_c \cdot \log \hat{y}_c = - \log \hat{y}_c$$

$$= - \log \frac{\exp(z_c)}{\sum_{j=1}^K \exp(z_j)} \leftarrow \text{softmax}$$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$



# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$

$$h(x) = x^2 + 1 \quad \text{and} \quad g(x) = x^2$$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$

$$h(x) = x^2 + 1 \quad \text{and} \quad g(x) = x^2$$

$$\frac{df}{dx} = 2(x^2 + 1) \cdot 2x$$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$

$$h(x) = x^2 + 1 \quad \text{and} \quad g(x) = x^2$$

$$\frac{df}{dx} = 2(x^2 + 1) \cdot 2x$$

$$f(x) = \ln(ax) = g(h(x))$$

$$h(x) = ax \quad \text{and} \quad g(x) = \ln(x)$$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$

$$h(x) = x^2 + 1 \quad \text{and} \quad g(x) = x^2$$

$$\frac{df}{dx} = 2(x^2 + 1) \cdot 2x$$

$$f(x) = \ln(ax) = g(h(x))$$

$$h(x) = ax \quad \text{and} \quad g(x) = \ln(x)$$

$$\frac{df}{dx} = \frac{1}{ax} \cdot a = \frac{1}{x}$$

# Backpropagation uses the chain rule

$$f(x) = g(h(x))$$

Chain rule:  $\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx}$

$$f(x) = (x^2 + 1)^2 = g(h(x)) \quad ??$$

$$h(x) = x^2 + 1 \quad \text{and} \quad g(x) = x^2$$

$$\frac{df}{dx} = 2(x^2 + 1) \cdot 2x$$

$$f(x) = \ln(ax) = g(h(x))$$

$$h(x) = ax \quad \text{and} \quad g(x) = \ln(x)$$

$$\frac{df}{dx} = \frac{1}{ax} \cdot a = \frac{1}{x}$$

$$\ln(ax) = \ln(a) + \ln(x)$$

# Multidimensional chain rule

$$\mathbf{x} \in \mathbb{R}^\ell$$

$$\mathbf{a} = h(\mathbf{x}), \quad \mathbb{R}^\ell \rightarrow \mathbb{R}^n$$


$$\mathbf{b} = g(\mathbf{a}), \quad \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\mathbb{R}^{\ell \times m} \quad \mathbb{R}^{\ell \times n} \quad \mathbb{R}^{n \times m}$$
$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \boxed{\frac{\partial \mathbf{a}}{\partial \mathbf{x}}} \cdot \boxed{\frac{\partial \mathbf{b}}{\partial \mathbf{a}}}$$

# Multidimensional chain rule

$$\begin{aligned}
 &\mathbf{x} \in \mathbb{R}^\ell \\
 &\mathbf{a} = h(\mathbf{x}), \quad \mathbb{R}^\ell \rightarrow \mathbb{R}^n \\
 &\mathbf{b} = g(\mathbf{a}), \quad \mathbb{R}^n \rightarrow \mathbb{R}^m
 \end{aligned}$$
  

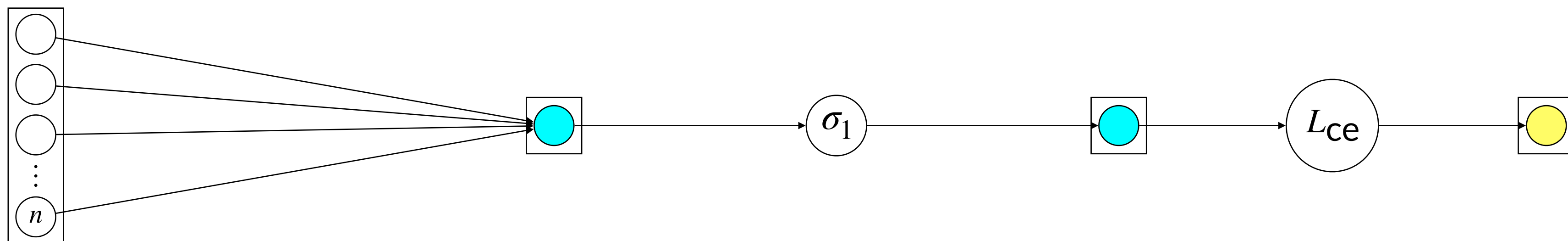
$$\mathbb{R}^{\ell \times m} \quad \mathbb{R}^{\ell \times n} \quad \mathbb{R}^{n \times m}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{x}} = \boxed{\frac{\partial \mathbf{a}}{\partial \mathbf{x}}} \cdot \boxed{\frac{\partial \mathbf{b}}{\partial \mathbf{a}}}$$


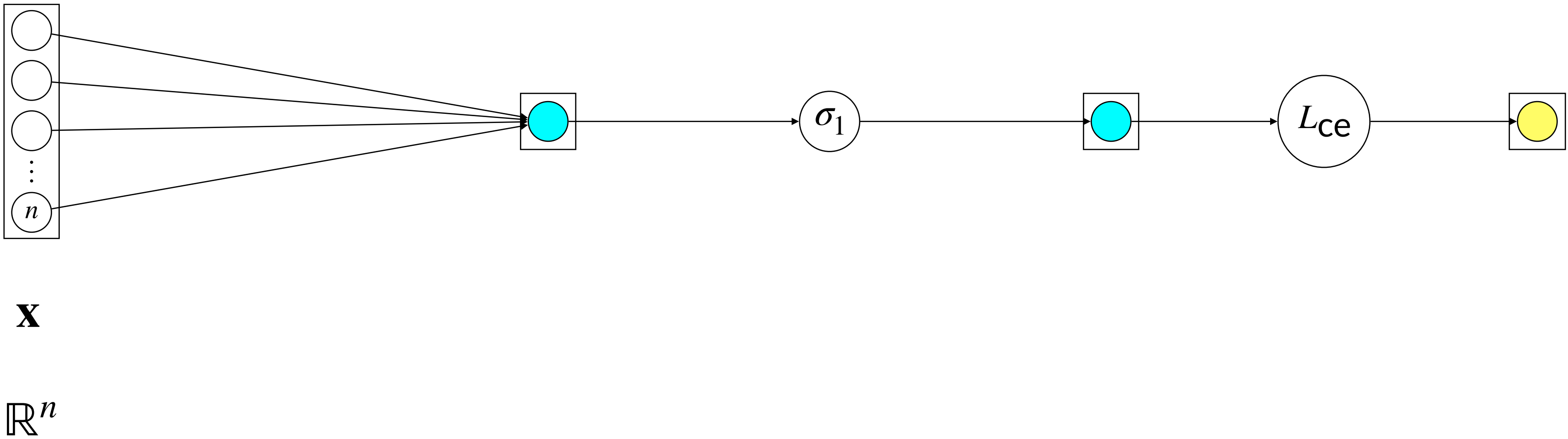
$$\begin{bmatrix}
 \frac{\partial b_1}{\partial a_1} & \frac{\partial b_2}{\partial a_1} & \frac{\partial b_3}{\partial a_1} & \cdots & \frac{\partial b_m}{\partial a_1} \\
 \frac{\partial b_1}{\partial a_2} & \frac{\partial b_2}{\partial a_2} & \frac{\partial b_3}{\partial a_2} & \cdots & \frac{\partial b_m}{\partial a_2} \\
 \frac{\partial b_1}{\partial a_3} & \frac{\partial b_2}{\partial a_3} & \frac{\partial b_3}{\partial a_3} & \cdots & \frac{\partial b_m}{\partial a_3} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 \frac{\partial b_1}{\partial a_n} & \frac{\partial b_2}{\partial a_n} & \frac{\partial b_3}{\partial a_n} & \cdots & \frac{\partial b_m}{\partial a_n}
 \end{bmatrix}$$



# Backprop and the chain rule in a 1-dimensional NN



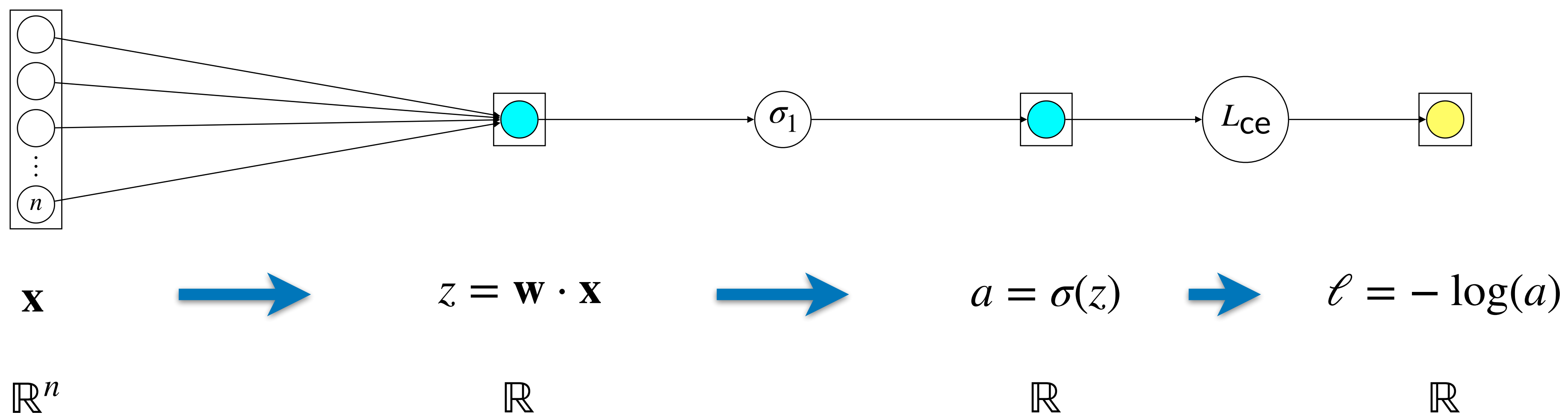
# Backprop and the chain rule in a 1-dimensional NN



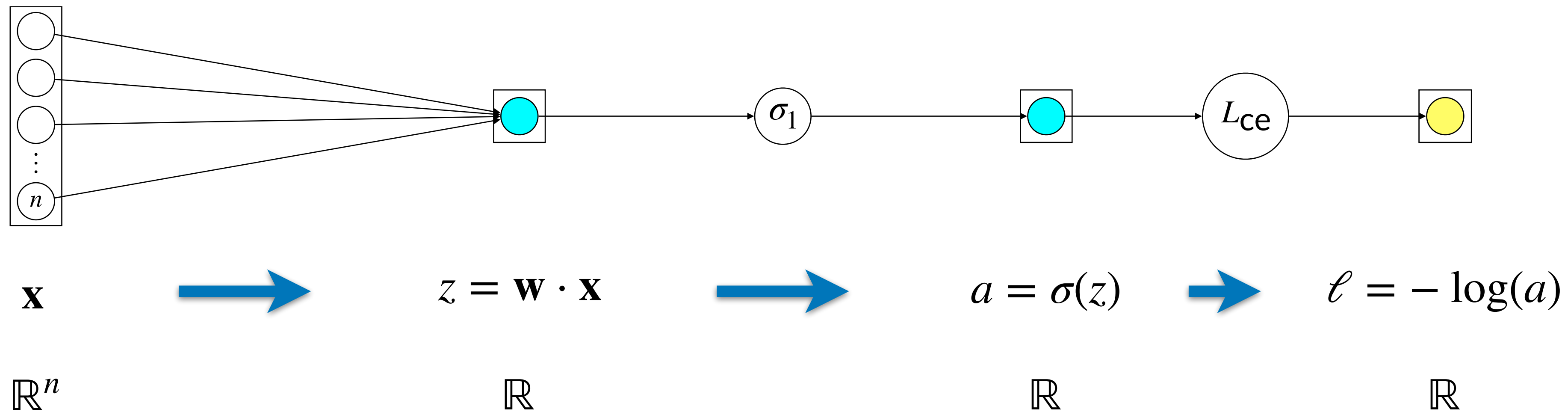




# Backprop and the chain rule in a 1-dimensional NN

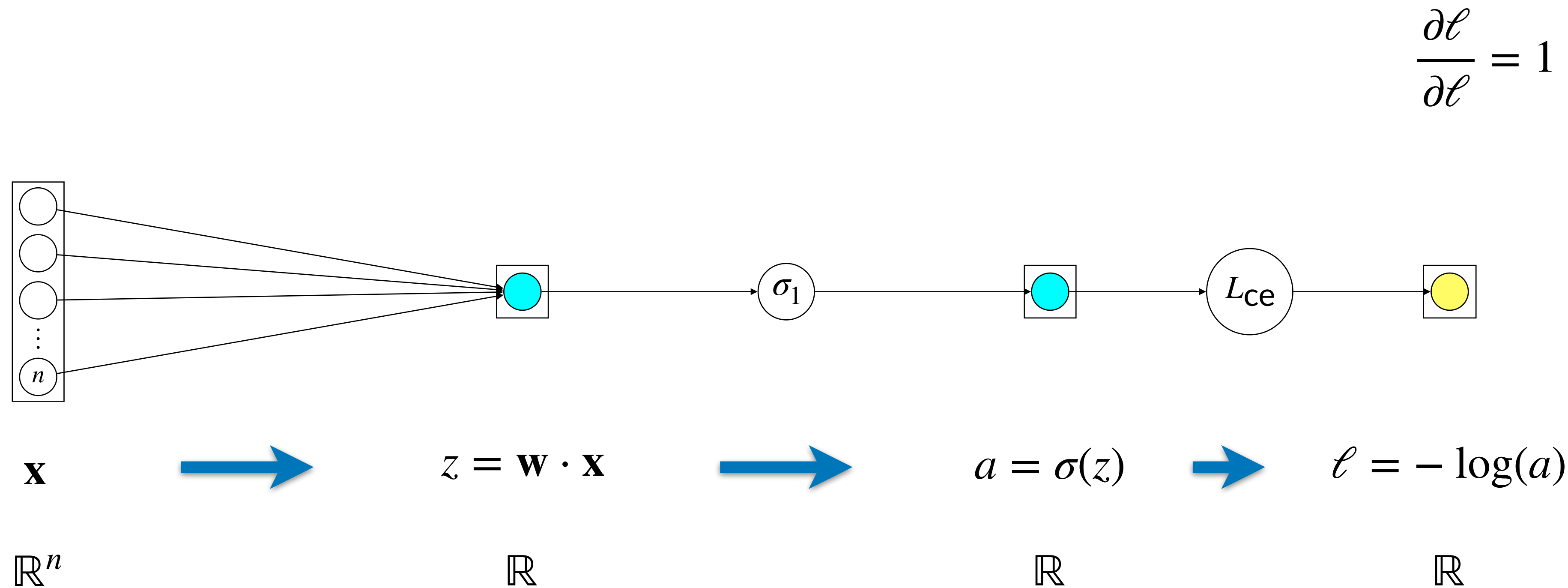


# Backprop and the chain rule in a 1-dimensional NN



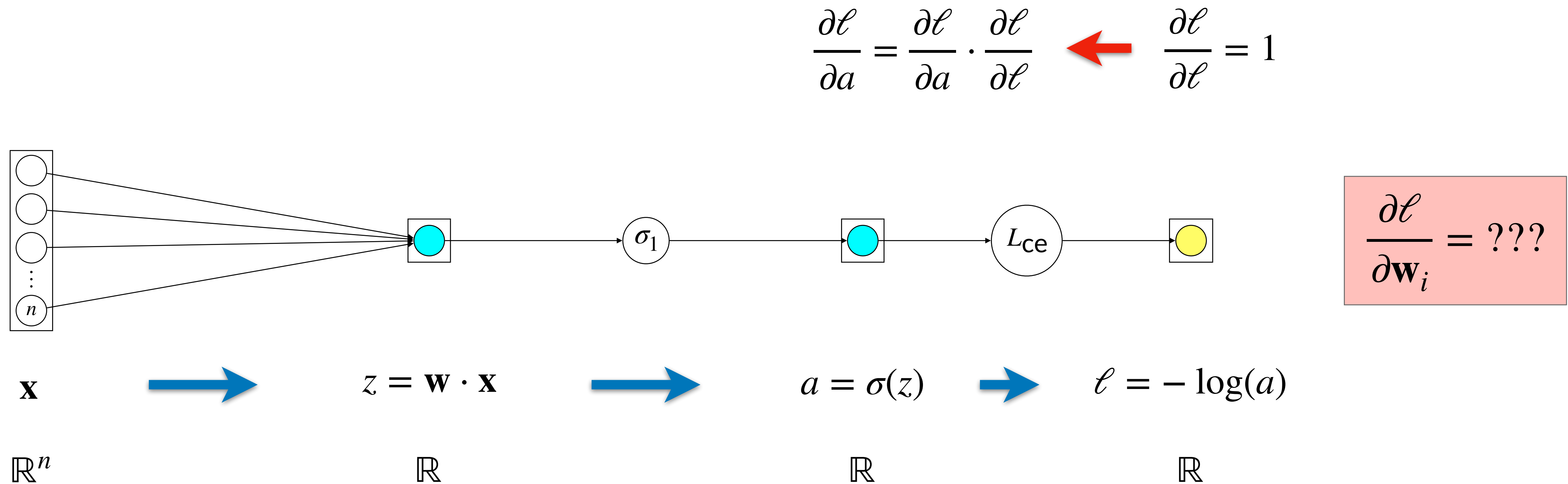
$$\frac{\partial \ell}{\partial \mathbf{w}_i} = ???$$

# Backprop and the chain rule in a 1-dimensional NN



$$\frac{\partial \ell}{\partial \mathbf{w}_i} = ???$$

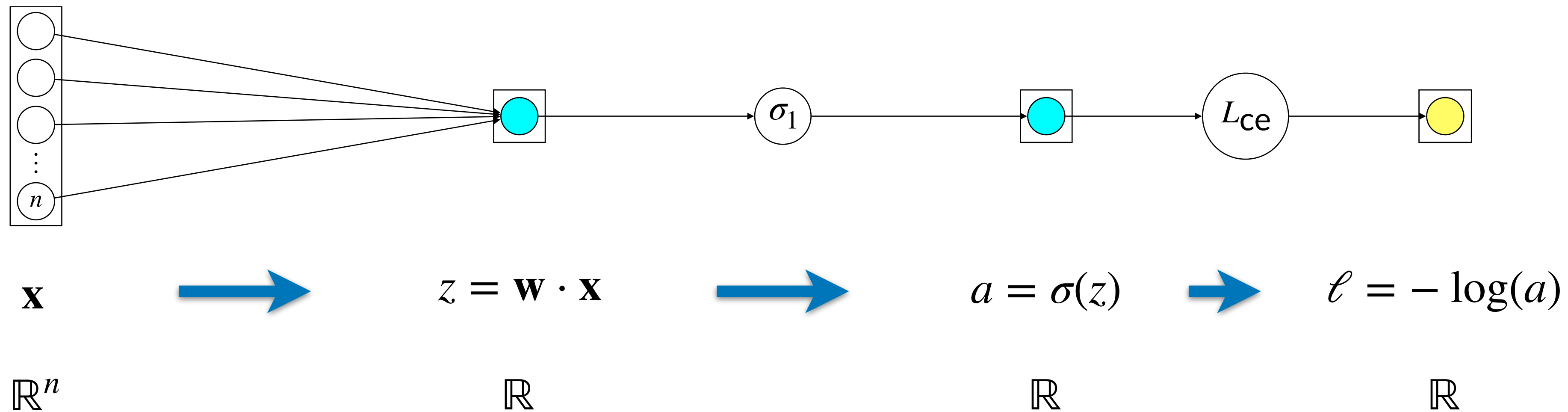
# Backprop and the chain rule in a 1-dimensional NN





# Backprop and the chain rule in a 1-dimensional NN

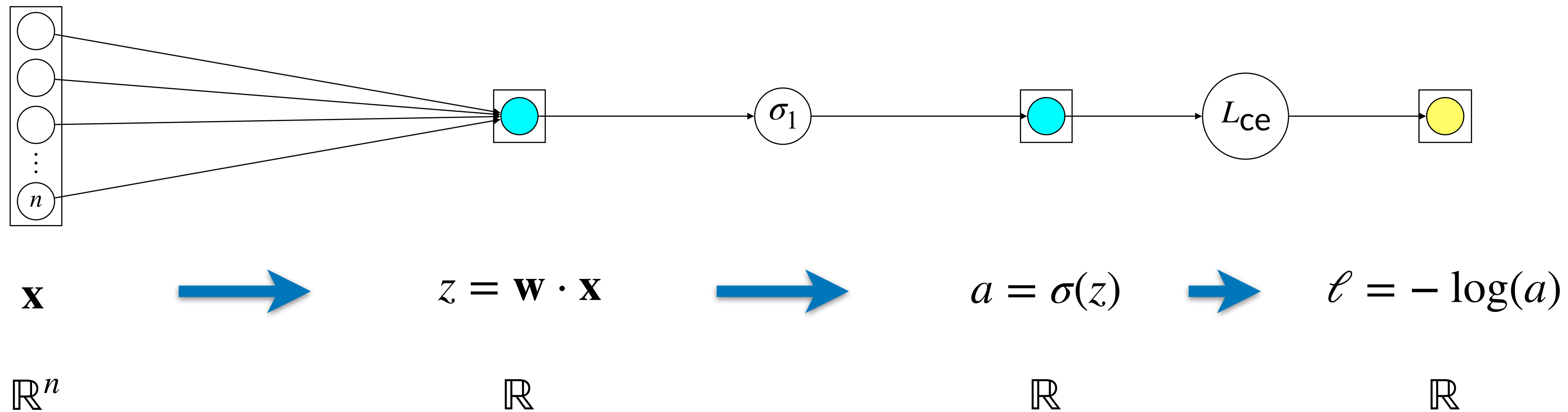
$$\frac{\partial \ell}{\partial z} = \frac{\partial a}{\partial z} \cdot \frac{\partial \ell}{\partial a} \quad \leftarrow \quad \frac{\partial \ell}{\partial a} = \frac{\partial \ell}{\partial a} \cdot \frac{\partial \ell}{\partial \ell} \quad \leftarrow \quad \frac{\partial \ell}{\partial \ell} = 1$$



$$\frac{\partial \ell}{\partial \mathbf{w}_i} = ???$$

# Backprop and the chain rule in a 1-dimensional NN

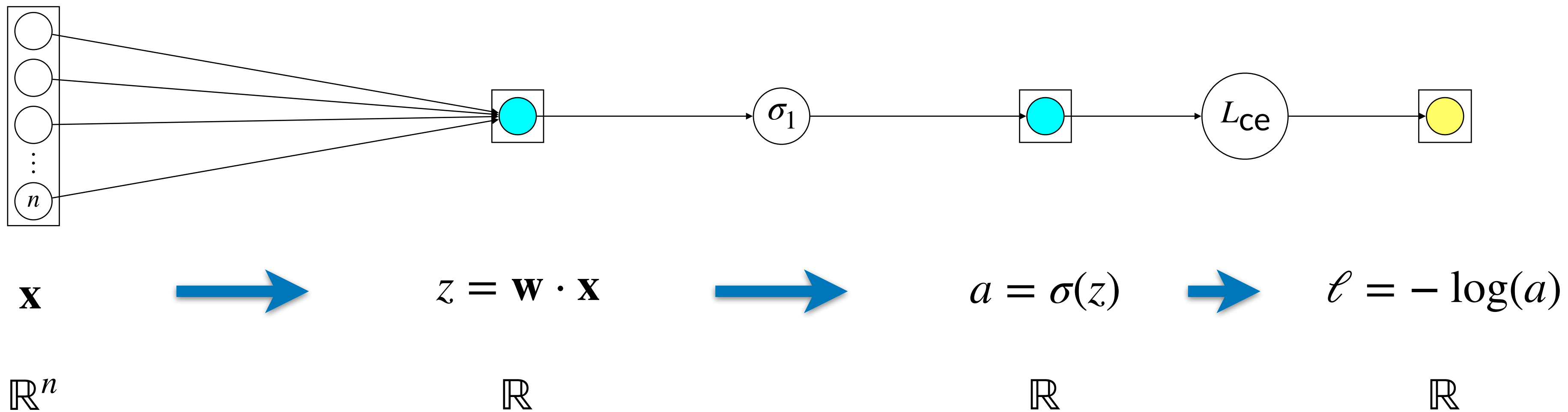
$$\frac{\partial \ell}{\partial \mathbf{w}_i} = \frac{\partial z}{\partial w_i} \cdot \frac{\partial \ell}{\partial z} \quad \leftarrow \quad \frac{\partial \ell}{\partial z} = \frac{\partial a}{\partial z} \cdot \frac{\partial \ell}{\partial a} \quad \leftarrow \quad \frac{\partial \ell}{\partial a} = \frac{\partial \ell}{\partial a} \cdot \frac{\partial \ell}{\partial \ell} \quad \leftarrow \quad \frac{\partial \ell}{\partial \ell} = 1$$



$$\frac{\partial \ell}{\partial \mathbf{w}_i} = ???$$

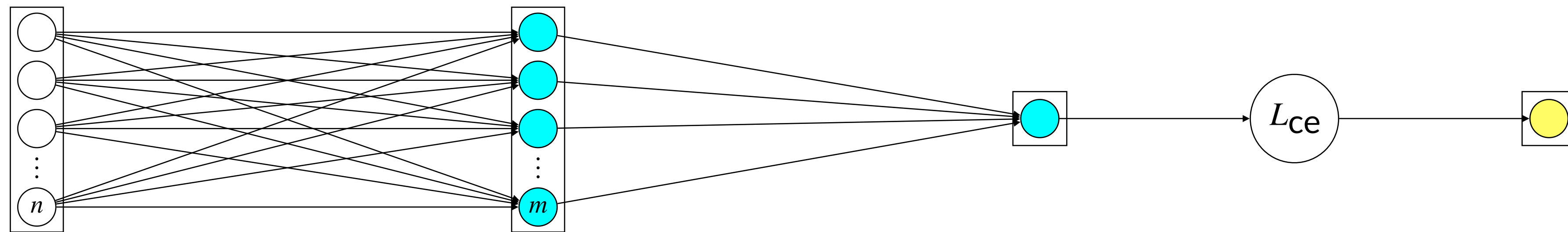
# Backprop and the chain rule in a 1-dimensional NN

$$\frac{\partial \ell}{\partial \mathbf{w}_i} = \frac{\partial z}{\partial w_i} \cdot \frac{\partial \ell}{\partial z} \leftarrow \frac{\partial \ell}{\partial z} = \frac{\partial a}{\partial z} \cdot \frac{\partial \ell}{\partial a} \leftarrow \frac{\partial \ell}{\partial a} = \frac{\partial \ell}{\partial a} \cdot \frac{\partial \ell}{\partial \ell} \leftarrow \frac{\partial \ell}{\partial \ell} = 1$$

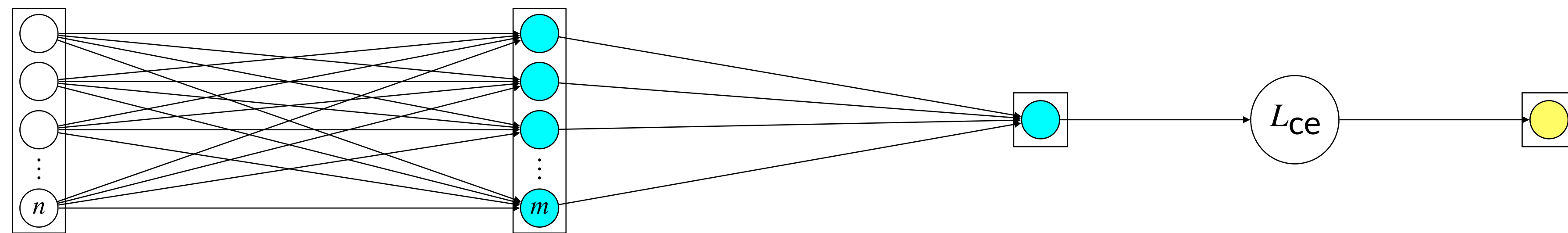


$$\frac{\partial \ell}{\partial \mathbf{w}_i} = ???$$

# Backprop and the chain rule in multiple dimensions



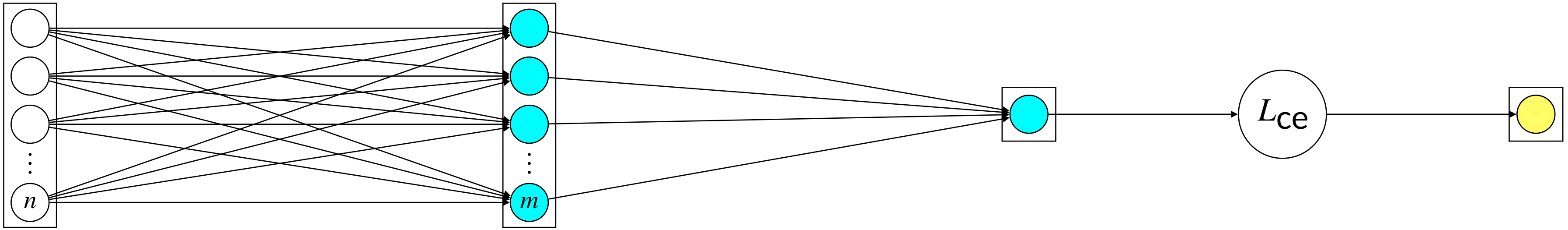
# Backprop and the chain rule in multiple dimensions



$\mathbf{a}^{[0]}$

$\mathbb{R}^n$

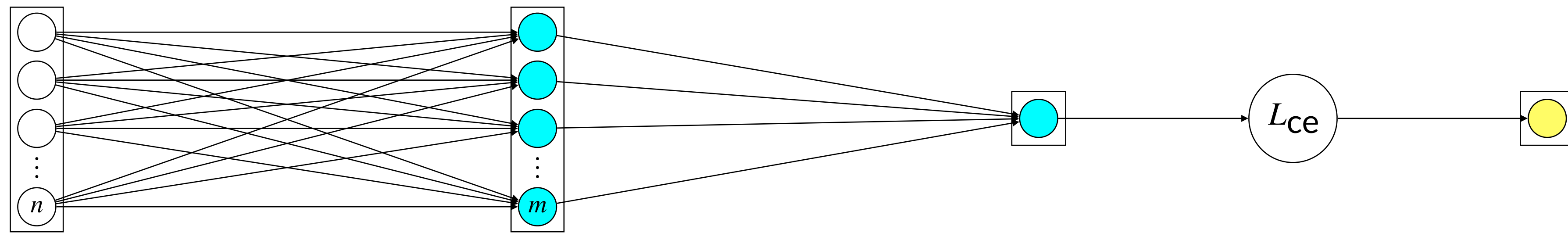
# Backprop and the chain rule in multiple dimensions



$\mathbf{a}^{[0]} \quad \longrightarrow \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]}$   
 $\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$

$\mathbb{R}^n$                        $\mathbb{R}^m$

# Backprop and the chain rule in multiple dimensions



$$\mathbf{a}^{[0]} \quad \longrightarrow \quad \begin{aligned} \mathbf{z}^{[1]} &= \mathbf{W}^{[1]} \mathbf{a}^{[0]} \\ \mathbf{a}^{[1]} &= \sigma(\mathbf{z}^{[1]}) \end{aligned} \quad \longrightarrow \quad \begin{aligned} z^{[2]} &= \mathbf{w}^{[2]} \mathbf{a}^{[1]} \\ a^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

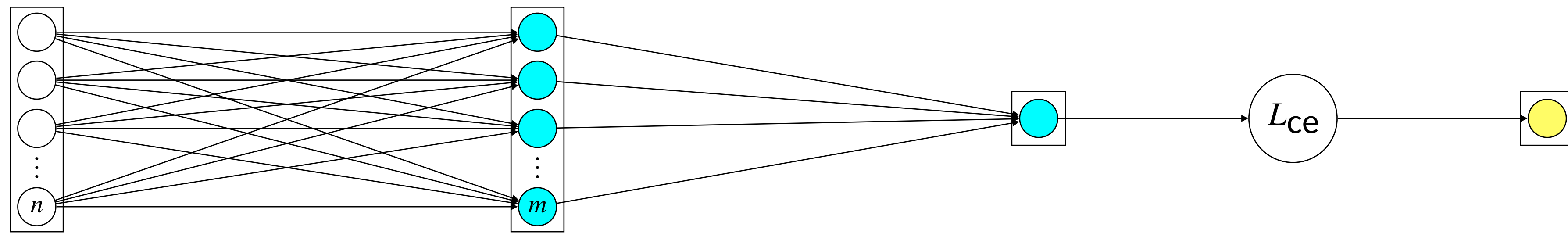
$\mathbb{R}^n$

$\mathbb{R}^m$

$\mathbb{R}$



# Backprop and the chain rule in multiple dimensions



$$\mathbf{a}^{[0]} \quad \longrightarrow \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} \quad \longrightarrow \quad z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} \quad \longrightarrow \quad \ell = -\log(a^{[2]})$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

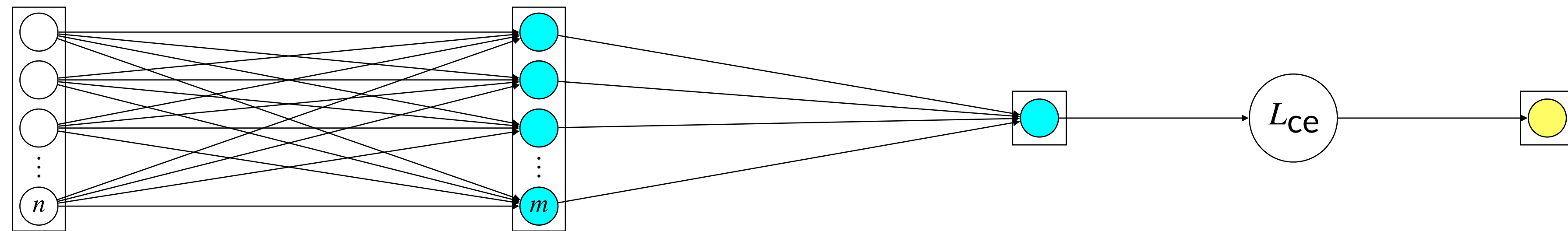
$\mathbb{R}^n$

$\mathbb{R}^m$

$\mathbb{R}$

$\mathbb{R}$

# Backprop and the chain rule in multiple dimensions



$$\mathbf{a}^{[0]} \quad \longrightarrow \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} \quad \longrightarrow \quad z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} \quad \longrightarrow \quad \ell = -\log(a^{[2]})$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

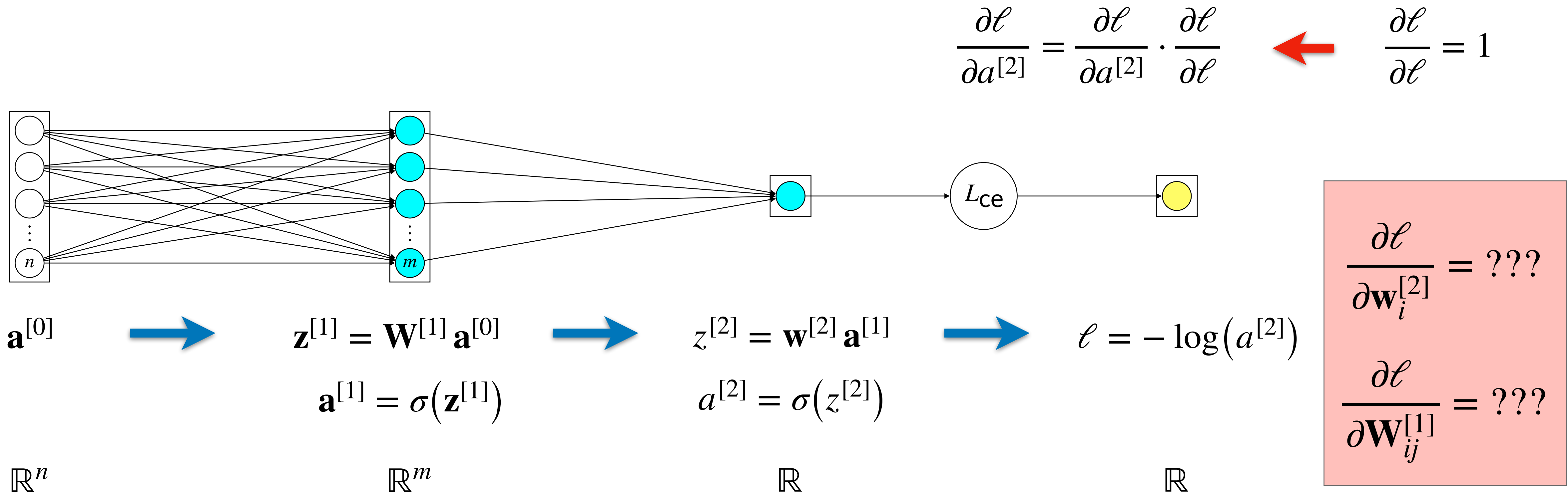
$$a^{[2]} = \sigma(z^{[2]})$$

 $\mathbb{R}^n$ 
 $\mathbb{R}^m$ 
 $\mathbb{R}$ 
 $\mathbb{R}$ 

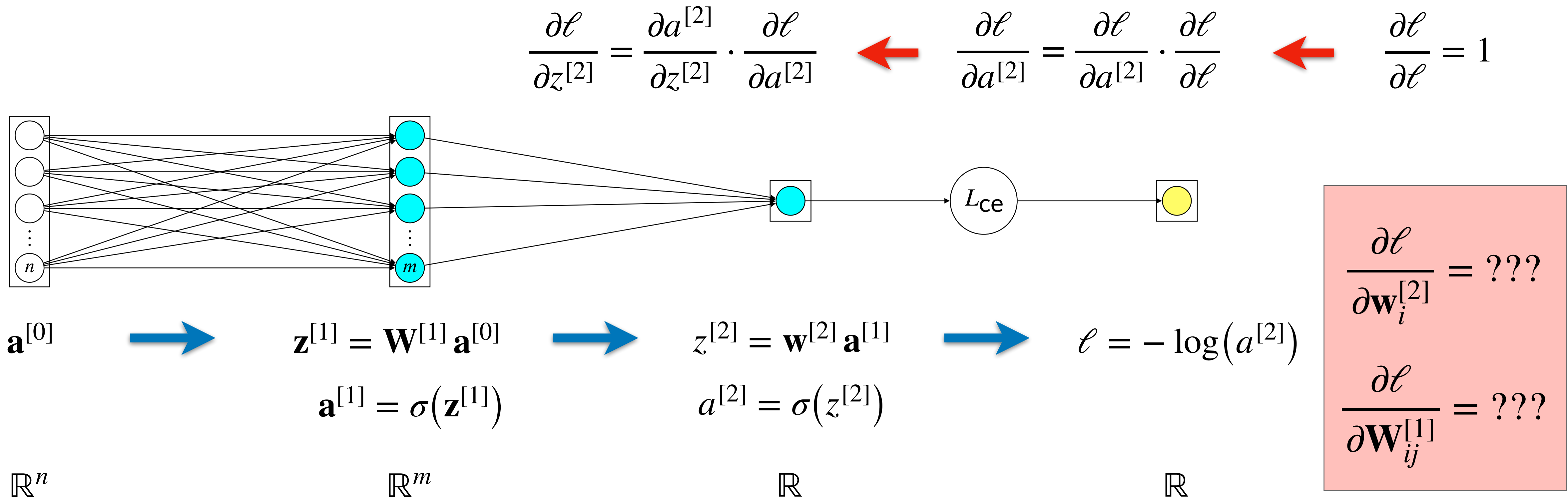
$$\frac{\partial \ell}{\partial \mathbf{w}_i^{[2]}} = ???$$

$$\frac{\partial \ell}{\partial \mathbf{W}_{ij}^{[1]}} = ???$$

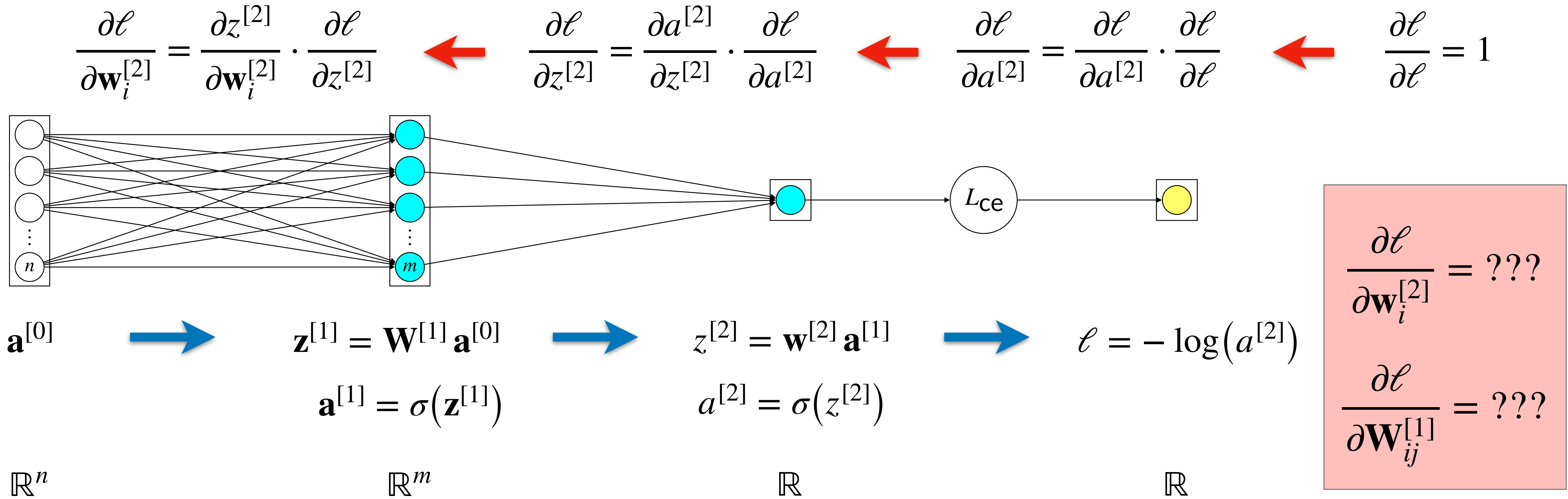
# Backprop and the chain rule in multiple dimensions



# Backprop and the chain rule in multiple dimensions



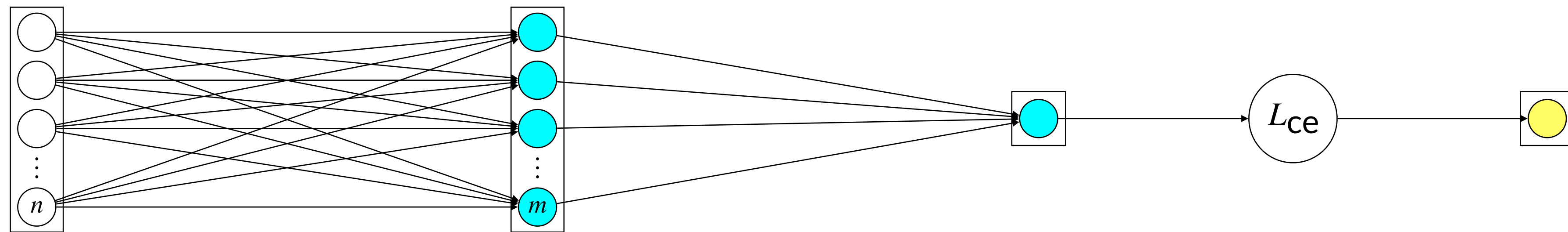
# Backprop and the chain rule in multiple dimensions



# Backprop and the chain rule in multiple dimensions

$$\frac{\partial \ell}{\partial \mathbf{W}_{ij}^{[1]}} = \frac{\partial \ell}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial \mathbf{a}_i^{[1]}} \cdot \frac{\partial \mathbf{a}_i^{[1]}}{\partial z_i^{[1]}} \cdot \frac{\partial z_i^{[1]}}{\partial \mathbf{W}_{ij}^{[1]}}$$

$$\frac{\partial \ell}{\partial \mathbf{w}_i^{[2]}} = \frac{\partial z^{[2]}}{\partial \mathbf{w}_i^{[2]}} \cdot \frac{\partial \ell}{\partial z^{[2]}} \quad \leftarrow \quad \frac{\partial \ell}{\partial z^{[2]}} = \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial \ell}{\partial a^{[2]}} \quad \leftarrow \quad \frac{\partial \ell}{\partial a^{[2]}} = \frac{\partial \ell}{\partial a^{[2]}} \cdot \frac{\partial \ell}{\partial \ell} \quad \leftarrow \quad \frac{\partial \ell}{\partial \ell} = 1$$



$$\mathbf{a}^{[0]} \quad \longrightarrow \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} \quad \longrightarrow \quad z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} \quad \longrightarrow \quad \ell = -\log(a^{[2]})$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

 $\mathbb{R}^n$ 
 $\mathbb{R}^m$ 
 $\mathbb{R}$ 
 $\mathbb{R}$ 

$$\frac{\partial \ell}{\partial \mathbf{w}_i^{[2]}} = ???$$

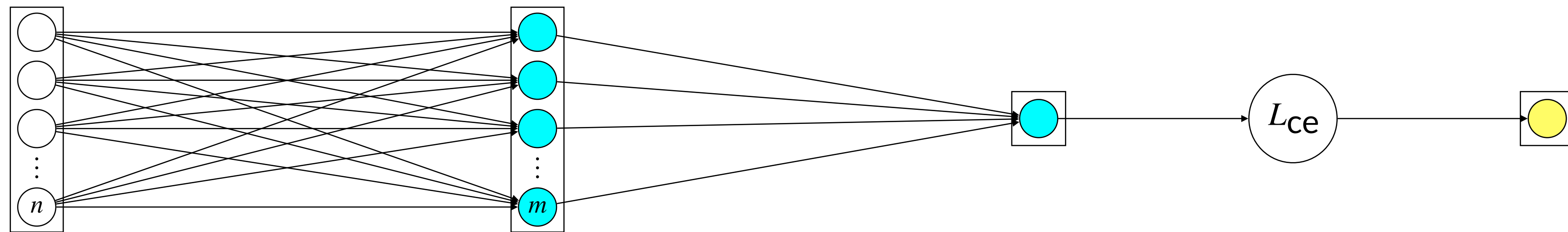
$$\frac{\partial \ell}{\partial \mathbf{W}_{ij}^{[1]}} = ???$$



# Backprop and the chain rule in multiple dimensions

$$\frac{\partial \ell}{\partial \mathbf{W}_{ij}^{[1]}} = \frac{\partial \ell}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial \mathbf{a}_i^{[1]}} \cdot \frac{\partial \mathbf{a}_i^{[1]}}{\partial z_i^{[1]}} \cdot \frac{\partial z_i^{[1]}}{\partial \mathbf{W}_{ij}^{[1]}}$$

$$\frac{\partial \ell}{\partial \mathbf{w}_i^{[2]}} = \frac{\partial z^{[2]}}{\partial \mathbf{w}_i^{[2]}} \cdot \frac{\partial \ell}{\partial z^{[2]}} \quad \leftarrow \quad \frac{\partial \ell}{\partial z^{[2]}} = \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial \ell}{\partial a^{[2]}} \quad \leftarrow \quad \frac{\partial \ell}{\partial a^{[2]}} = \frac{\partial \ell}{\partial a^{[2]}} \cdot \frac{\partial \ell}{\partial \ell} \quad \leftarrow \quad \frac{\partial \ell}{\partial \ell} = 1$$



$$\mathbf{a}^{[0]} \quad \longrightarrow \quad \mathbf{z}^{[1]} = \mathbf{W}^{[1]} \mathbf{a}^{[0]} \quad \longrightarrow \quad z^{[2]} = \mathbf{w}^{[2]} \mathbf{a}^{[1]} \quad \longrightarrow \quad \ell = -\log(a^{[2]})$$

$$\mathbf{a}^{[1]} = \sigma(\mathbf{z}^{[1]})$$

$$a^{[2]} = \sigma(z^{[2]})$$

 $\mathbb{R}^n$ 
 $\mathbb{R}^m$ 
 $\mathbb{R}$ 
 $\mathbb{R}$ 

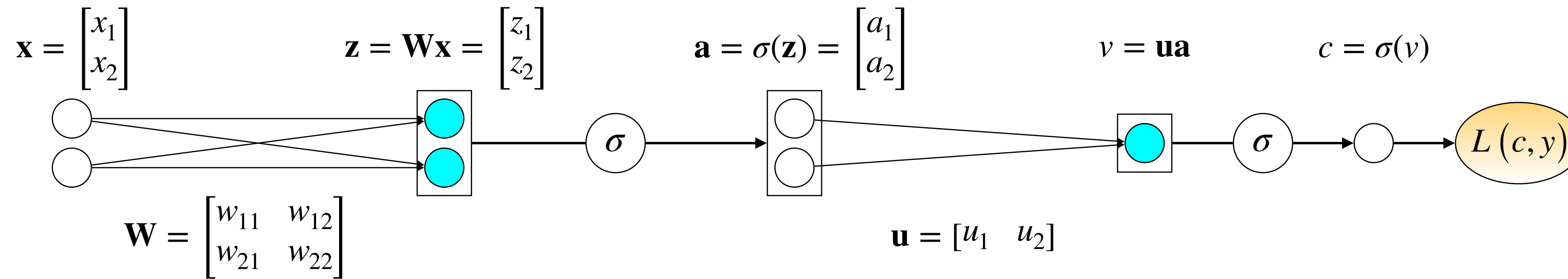
$$\frac{\partial \ell}{\partial \mathbf{w}_i^{[2]}} = ???$$

$$\frac{\partial \ell}{\partial \mathbf{W}_{ij}^{[1]}} = ???$$





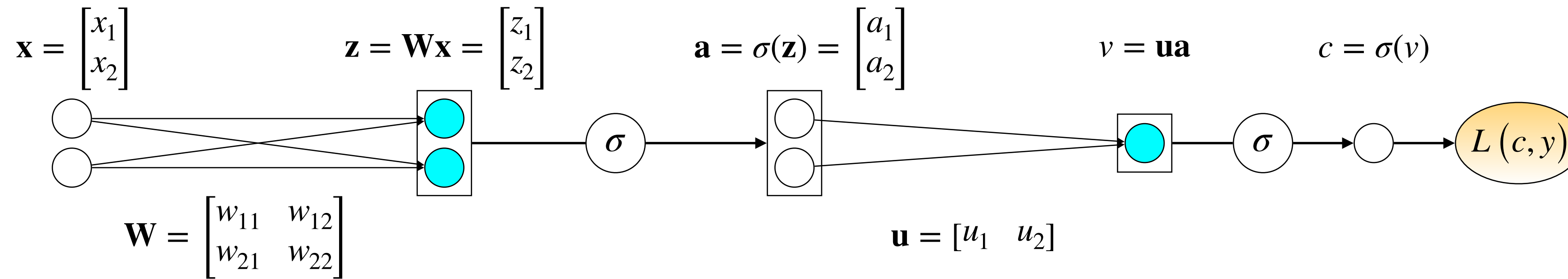
# Backprop (toy) example (1)



- ▶ Feedforward NN with 2 layers with 2 and 1 units
- ▶ Just 2 inputs and 1 output (*binary classification*)
- ▶ No bias terms and different letters used for different variables to simplify notation in upcoming slides
- ▶ Sigmoid (*logistic*) activation function everywhere

We want to update  $\mathbf{W}$  and  $\mathbf{u}$  with respect to the loss  $L(c, y)$

# Backprop (toy) example (2)



Cross-entropy loss  
for binary classification

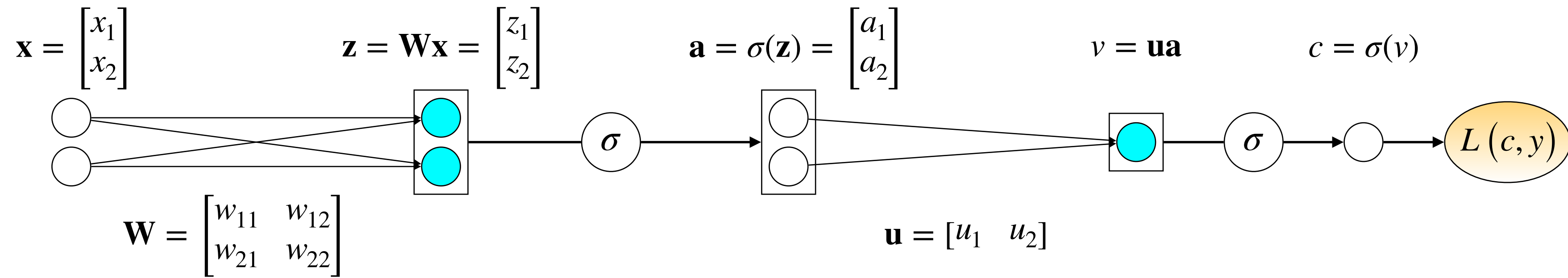
$$L(c, y) = - \left[ y \ln(c) + (1 - y) \ln(1 - c) \right]$$

Understanding how a matrix  
operation looks like is key in  
getting the derivatives right

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \mathbf{W} \cdot \mathbf{x} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$z_i = \sum_{j=1}^2 w_{ij} \cdot x_j$$

# Backprop (toy) example (3)

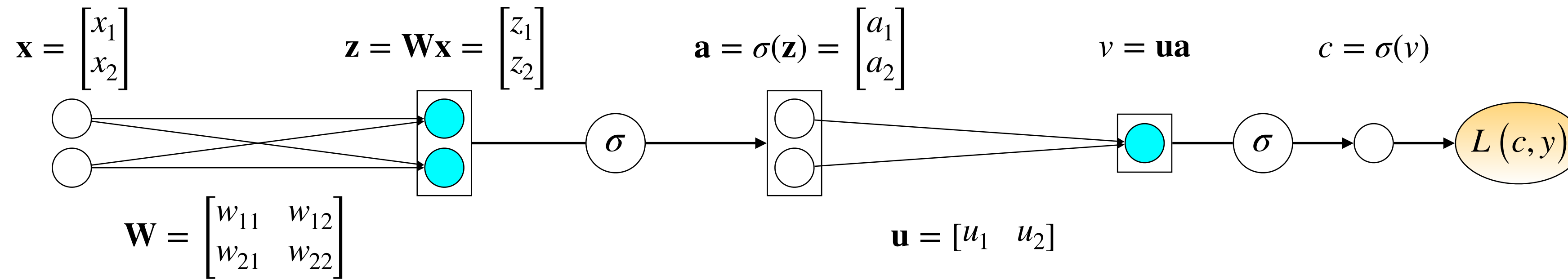


The derivative of the sigmoid activation is *neat*

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$\frac{d\sigma}{dx} = \frac{-\exp(-x)}{(1 + \exp(-x))^2}$$

# Backprop (toy) example (3)

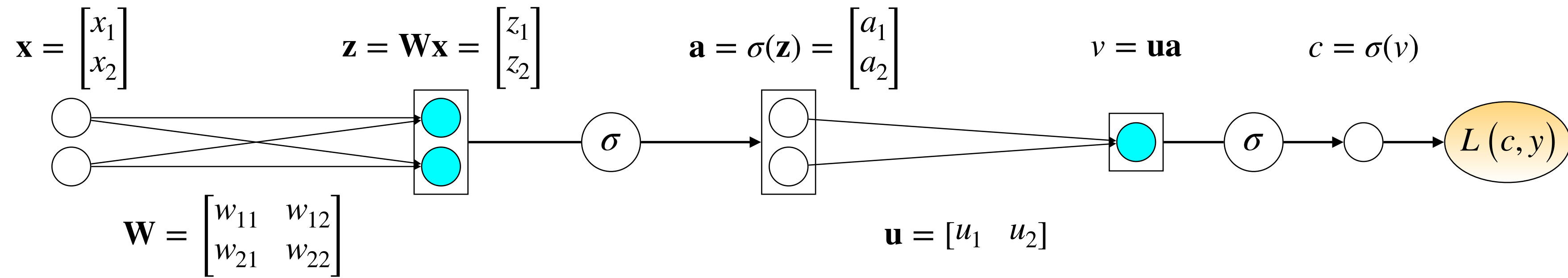


The derivative of the sigmoid activation is *neat*

$$\sigma(x) = \frac{1}{1 - \exp(-x)}$$

$$\frac{d\sigma}{dx} = \frac{-\exp(-x)}{(1 - \exp(-x))^2} = \frac{1}{1 - \exp(-x)} \cdot \frac{-\exp(-x)}{1 - \exp(-x)}$$

# Backprop (toy) example (3)



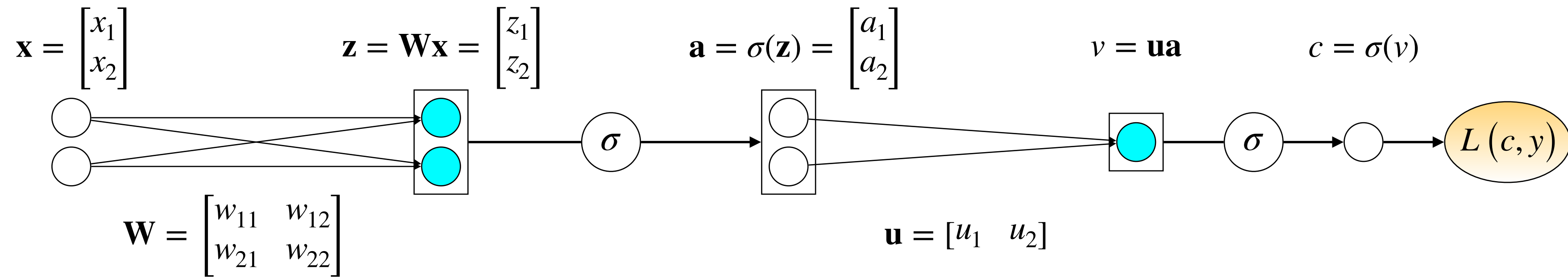
The derivative of the sigmoid activation is *neat*

$$\sigma(x) = \frac{1}{1 - \exp(-x)}$$

$$\frac{d\sigma}{dx} = \frac{-\exp(-x)}{(1 - \exp(-x))^2} = \frac{1}{1 - \exp(-x)} \cdot \frac{-\exp(-x)}{1 - \exp(-x)} = \frac{1}{1 - \exp(-x)} \cdot \frac{\mathbf{1 - \exp(-x) - 1}}{1 - \exp(-x)}$$

add and subtract 1

# Backprop (toy) example (3)



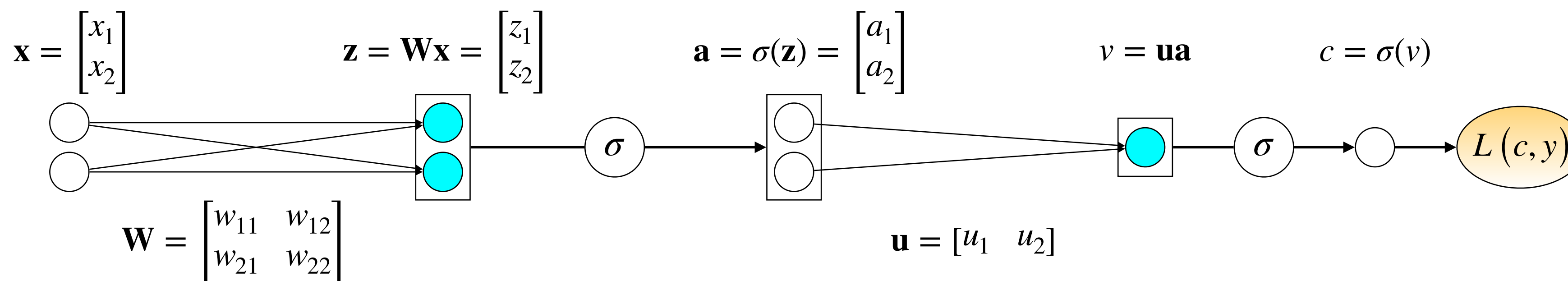
The derivative of the sigmoid activation is *neat*

$$\sigma(x) = \frac{1}{1 - \exp(-x)}$$

add and subtract 1

$$\begin{aligned}
 \frac{d\sigma}{dx} &= \frac{-\exp(-x)}{(1 - \exp(-x))^2} = \frac{1}{1 - \exp(-x)} \cdot \frac{-\exp(-x)}{1 - \exp(-x)} = \frac{1}{1 - \exp(-x)} \cdot \frac{\mathbf{1} - \exp(-\mathbf{x}) - \mathbf{1}}{1 - \exp(-x)} \\
 &= \frac{1}{1 - \exp(-x)} \cdot \left( \frac{1 - \exp(-x)}{1 - \exp(-x)} - \frac{1}{1 - \exp(-x)} \right)
 \end{aligned}$$

# Backprop (toy) example (3)



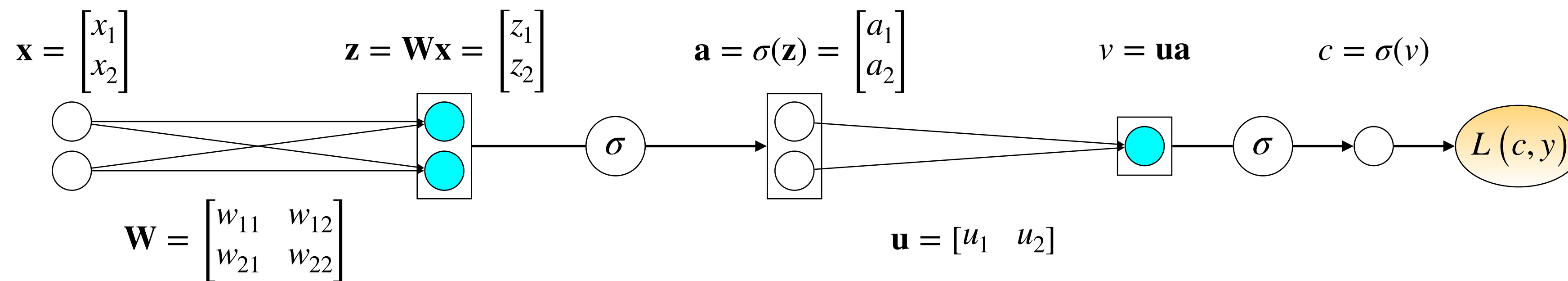
The derivative of the sigmoid activation is neat

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$\begin{aligned}
 \frac{d\sigma}{dx} &= \frac{-\exp(-x)}{(1 + \exp(-x))^2} \\
 &= \frac{1}{1 + \exp(-x)} \cdot \frac{-\exp(-x)}{1 + \exp(-x)} \\
 &= \frac{1}{1 + \exp(-x)} \cdot \frac{1 - \exp(-x) - 1}{1 + \exp(-x)} \quad \text{add and subtract 1} \\
 &= \frac{1}{1 + \exp(-x)} \cdot \left( \frac{1 - \exp(-x)}{1 + \exp(-x)} - \frac{1}{1 + \exp(-x)} \right) \\
 &= \sigma(x) \cdot (1 - \sigma(x))
 \end{aligned}$$



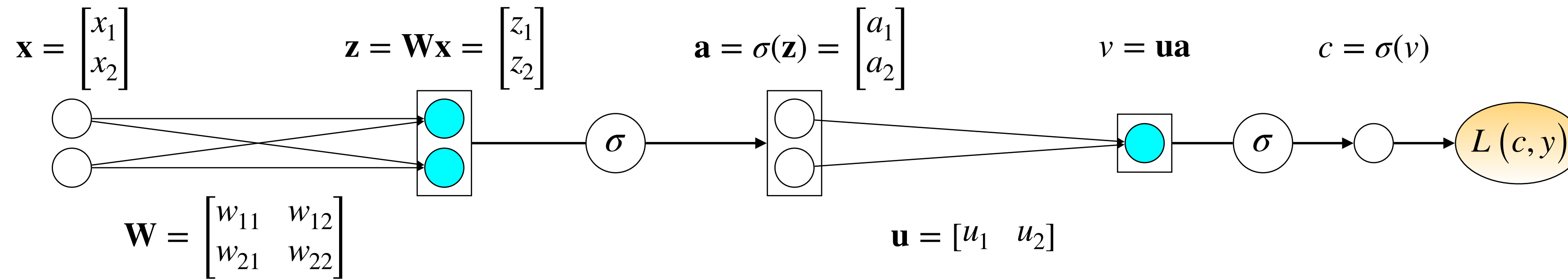
# Backprop (toy) example (4)



We first want to obtain this

$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial u_i}$$

# Backprop (toy) example (4)



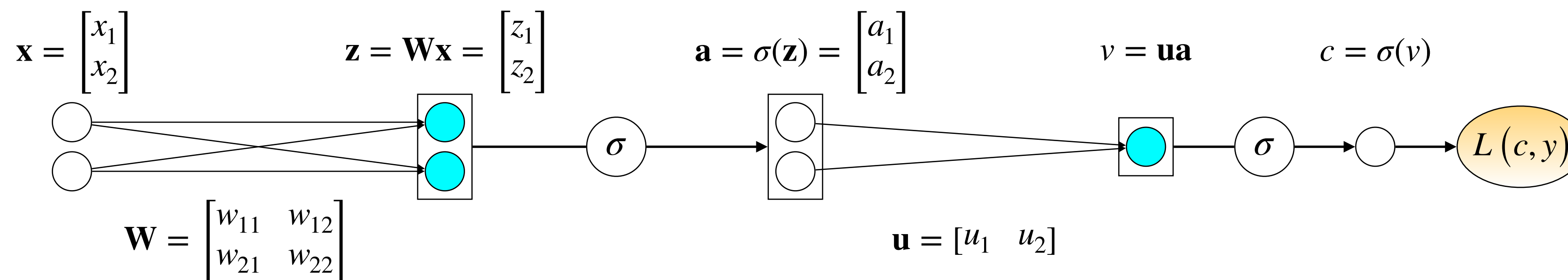
We first want to obtain this

$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial u_i}$$

Chain rule

$$= \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

# Backprop (toy) example (5)

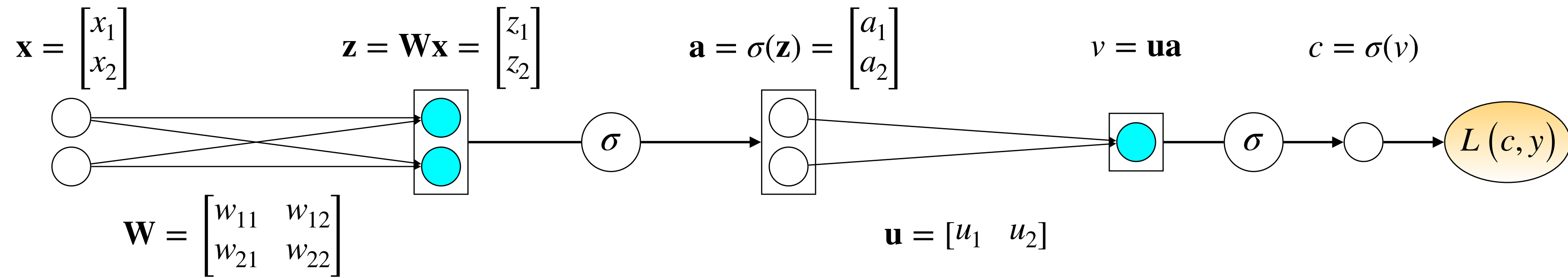


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$L(c, y) = -y \ln(c) - (1 - y) \ln(1 - c)$$

$$\frac{\partial L}{\partial c} =$$

# Backprop (toy) example (5)

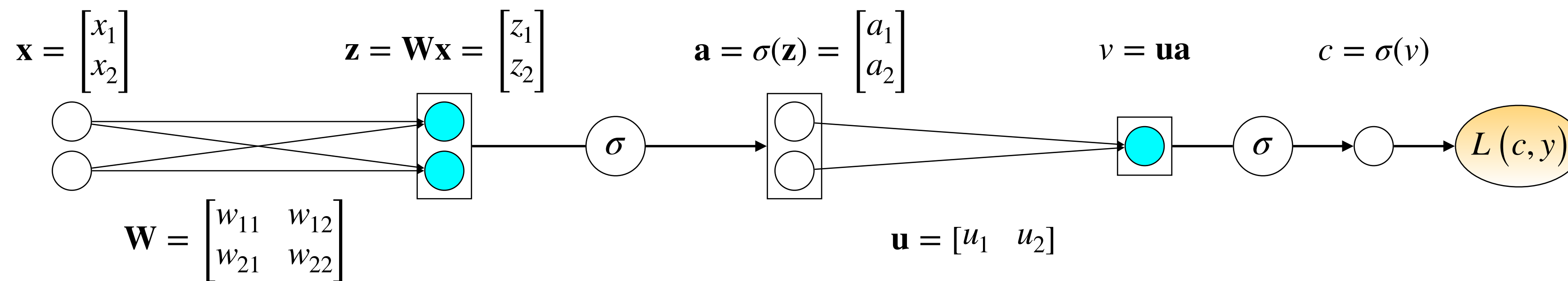


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$L(c, y) = -y \ln(c) - (1 - y) \ln(1 - c)$$

$$\frac{\partial L}{\partial c} = -y \cdot \frac{1}{c}$$

# Backprop (toy) example (5)

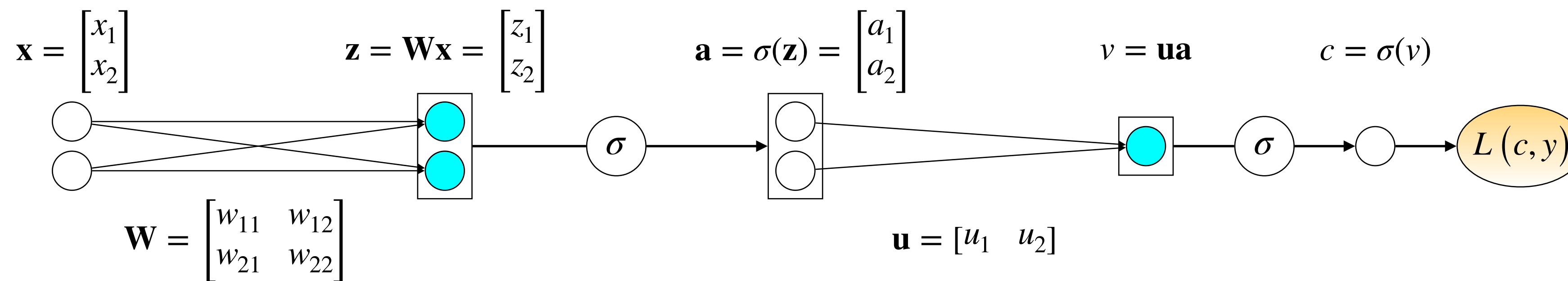


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$L(c, y) = -y \ln(c) - (1 - y) \ln(1 - c)$$

$$\frac{\partial L}{\partial c} = -y \cdot \frac{1}{c} - (1 - y) \cdot \frac{1}{1 - c} \cdot (-1)$$

# Backprop (toy) example (5)

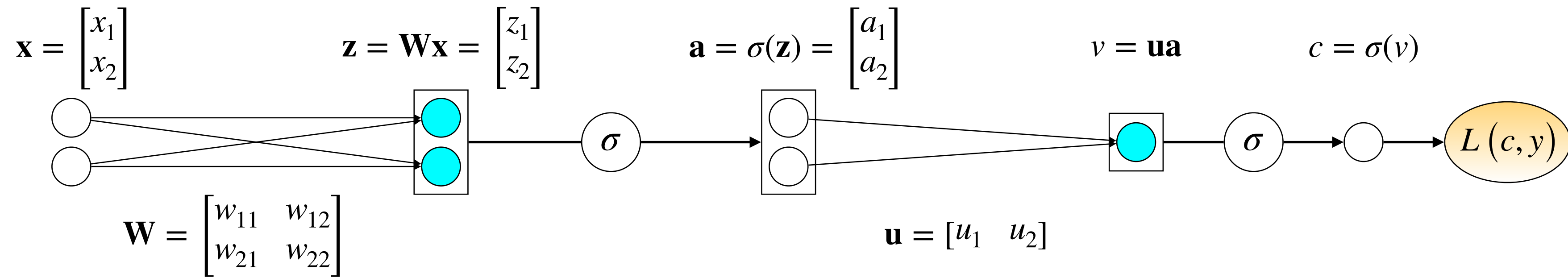


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$L(c, y) = -y \ln(c) - (1 - y) \ln(1 - c)$$

$$\frac{\partial L}{\partial c} = -y \cdot \frac{1}{c} - (1 - y) \cdot \frac{1}{1 - c} \cdot (-1) = \frac{c - y}{c(1 - c)}$$

# Backprop (toy) example (6)



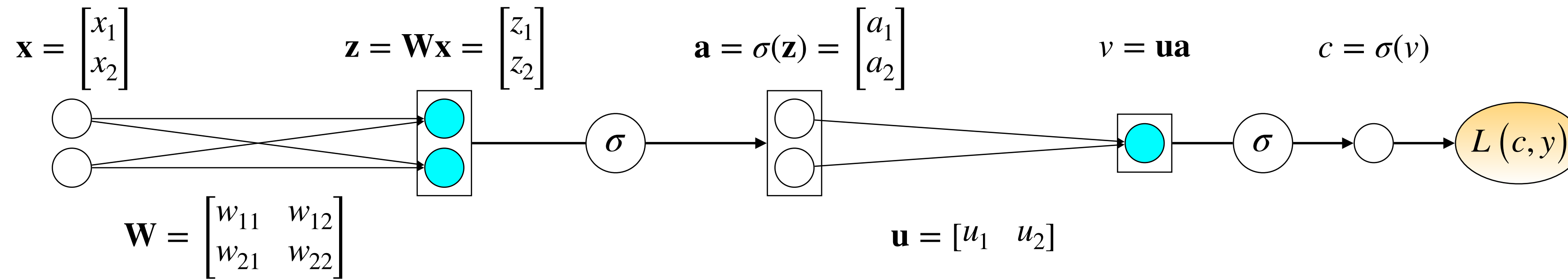
$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

reminder  $\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial c}{\partial v} =$$



# Backprop (toy) example (6)

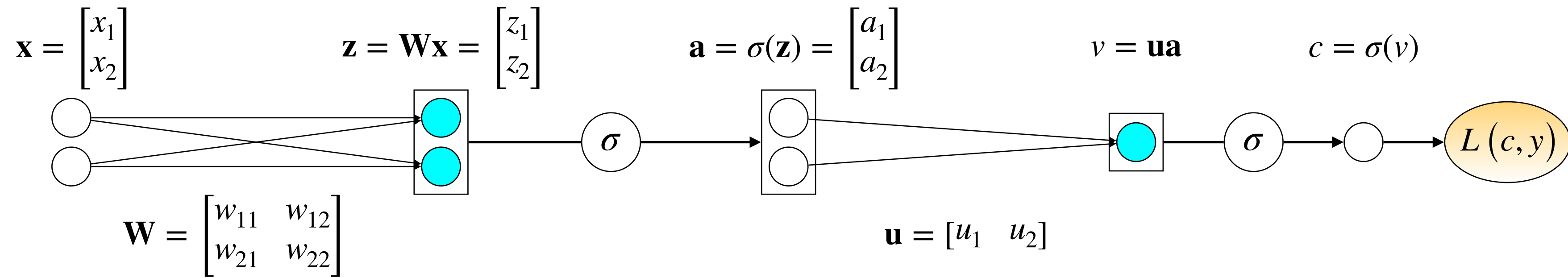


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

reminder  $\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial c}{\partial v} = \frac{\partial \sigma(v)}{\partial v}$$

# Backprop (toy) example (6)

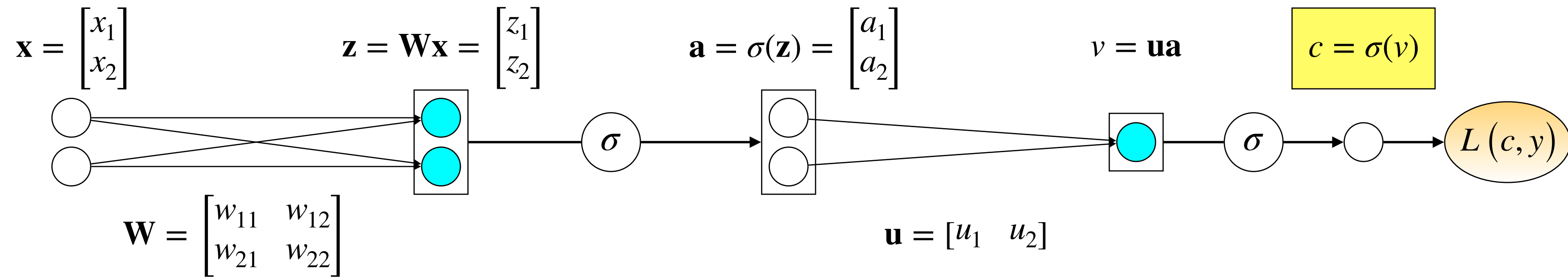


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

reminder  $\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial c}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v) \cdot (1 - \sigma(v))$$

# Backprop (toy) example (6)

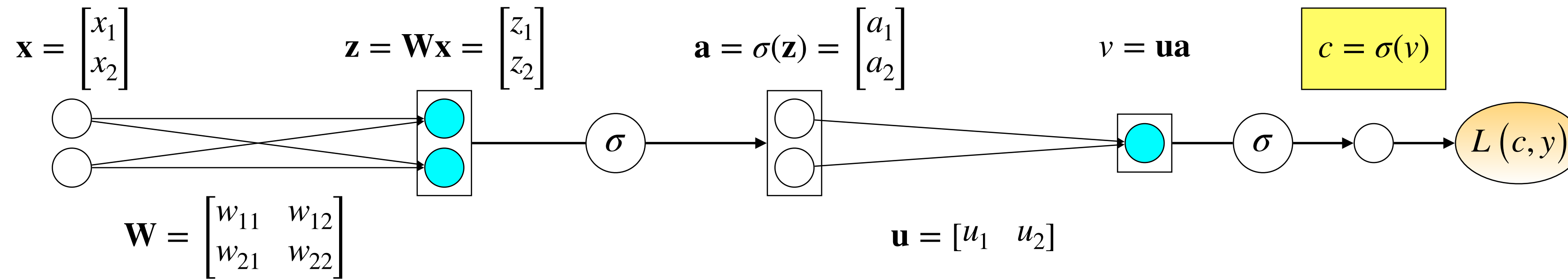


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

reminder  $\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial c}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v) \cdot (1 - \sigma(v))$$

# Backprop (toy) example (6)

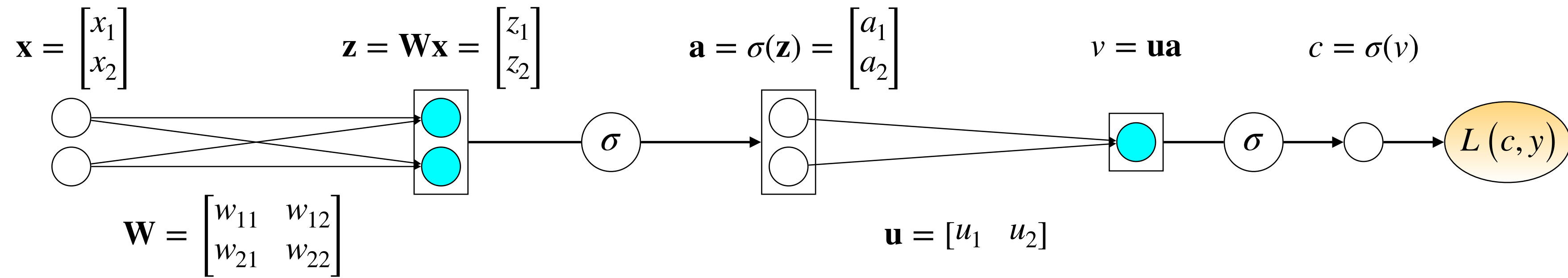


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

reminder  $\frac{d\sigma}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

$$\frac{\partial c}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v) \cdot (1 - \sigma(v)) = c \cdot (1 - c)$$

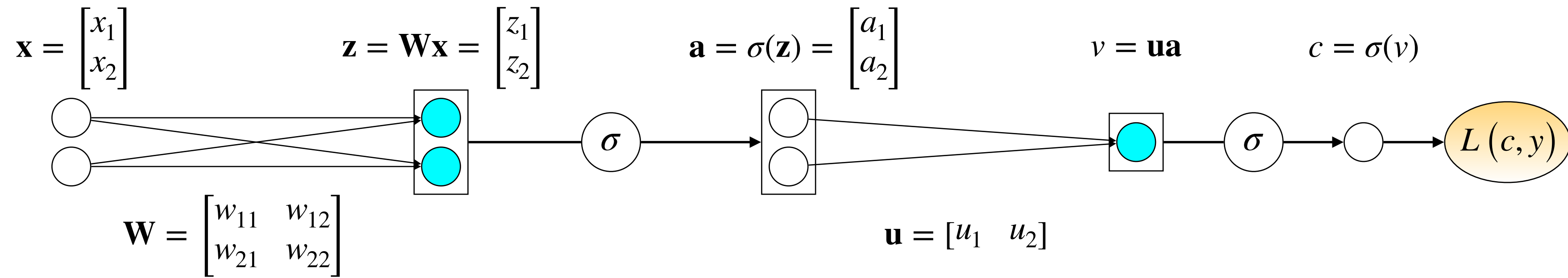
# Backprop (toy) example (7)



$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$\frac{\partial v}{\partial u_i} = \frac{\partial(\mathbf{u}\mathbf{a})}{\partial u_i}$$

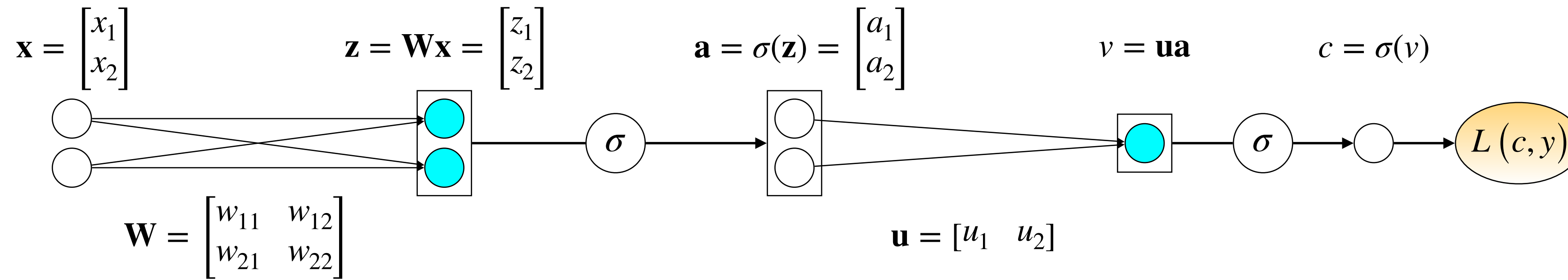
# Backprop (toy) example (7)



$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$\frac{\partial v}{\partial u_i} = \frac{\partial(\mathbf{u}\mathbf{a})}{\partial u_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial u_i}$$

# Backprop (toy) example (7)

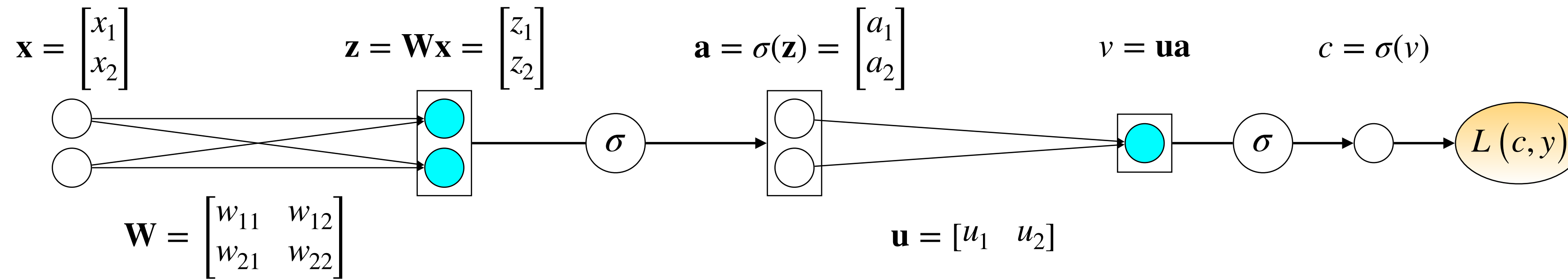


$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$\frac{\partial v}{\partial u_i} = \frac{\partial(\mathbf{u}\mathbf{a})}{\partial u_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial u_i} = a_i$$



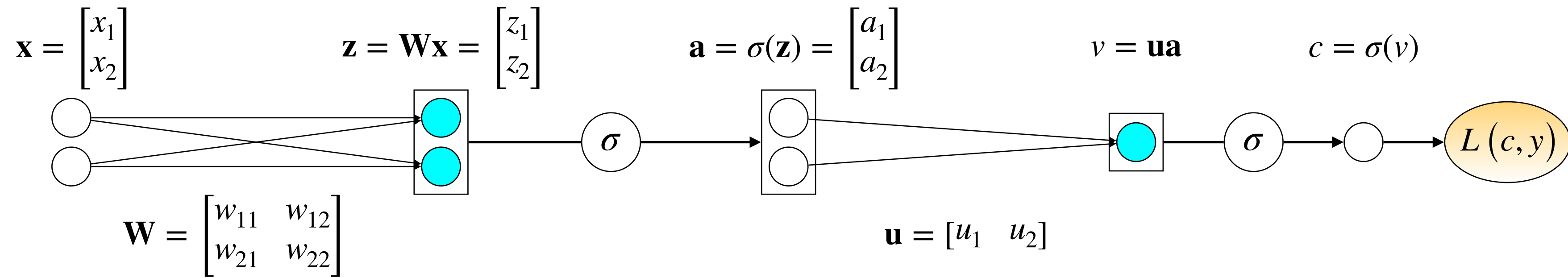
# Backprop (toy) example (8)



$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i}$$

$$\frac{\partial L}{\partial c} = \frac{c - y}{c \cdot (1 - c)} \quad \frac{\partial c}{\partial v} = c \cdot (1 - c) \quad \frac{\partial v}{\partial u_i} = a_i$$

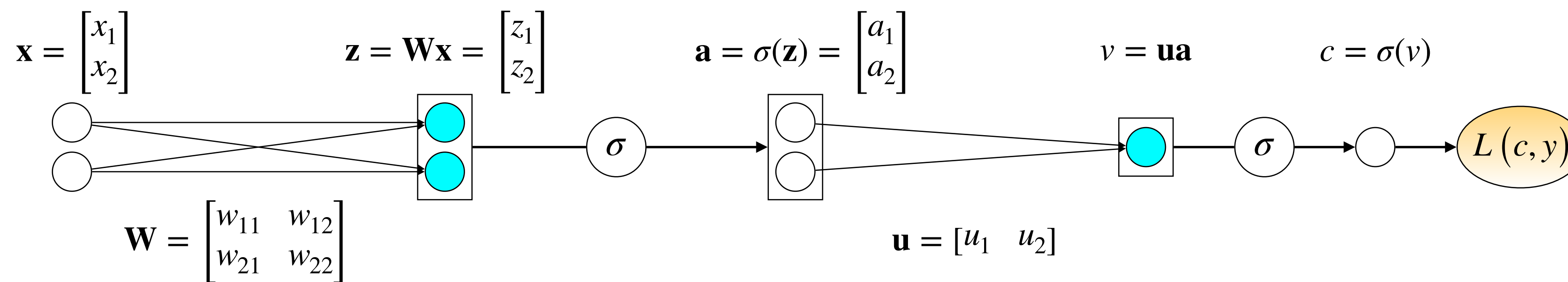
# Backprop (toy) example (8)



$$\frac{\partial L}{\partial u_i} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial u_i} = (c - y) \cdot a_i$$

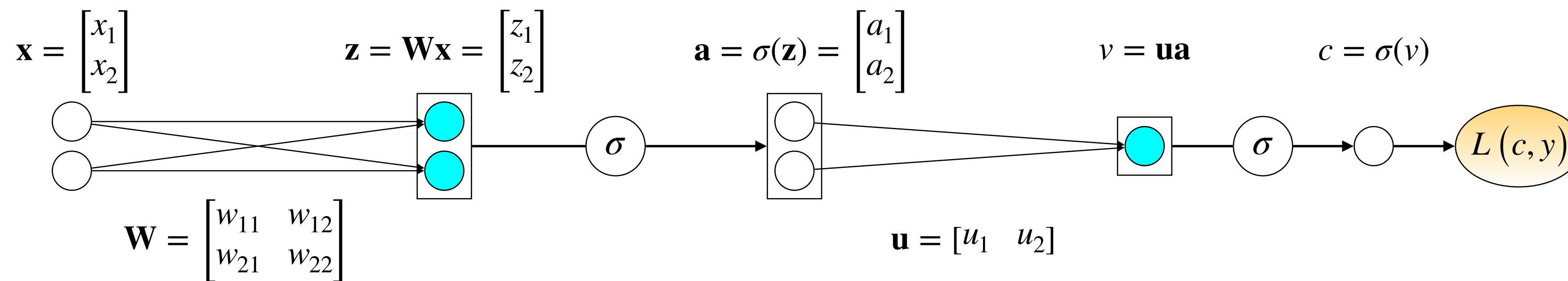
$$\frac{\partial L}{\partial c} = \frac{c - y}{c \cdot (1 - c)} \quad \frac{\partial c}{\partial v} = c \cdot (1 - c) \quad \frac{\partial v}{\partial u_i} = a_i$$

# Backprop (toy) example (9)



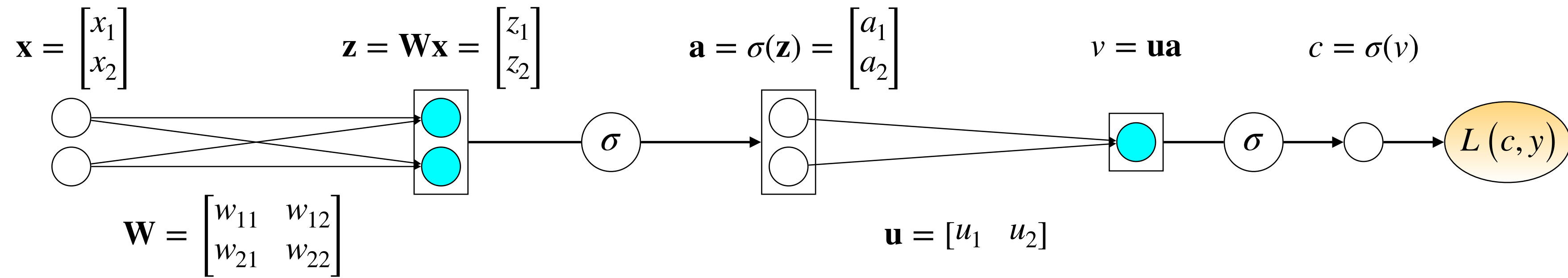
$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

# Backprop (toy) example (9)



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

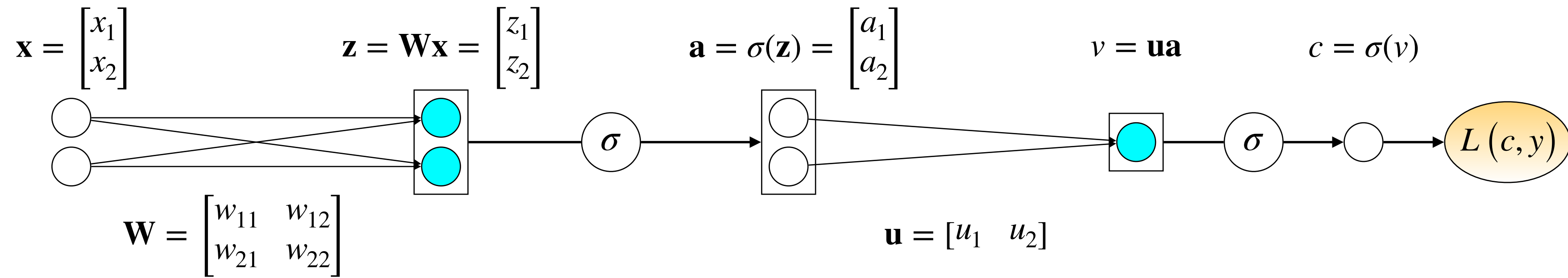
# Backprop (toy) example (9)



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

$$\frac{\partial v}{\partial a_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial a_i} = u_i$$

# Backprop (toy) example (9)

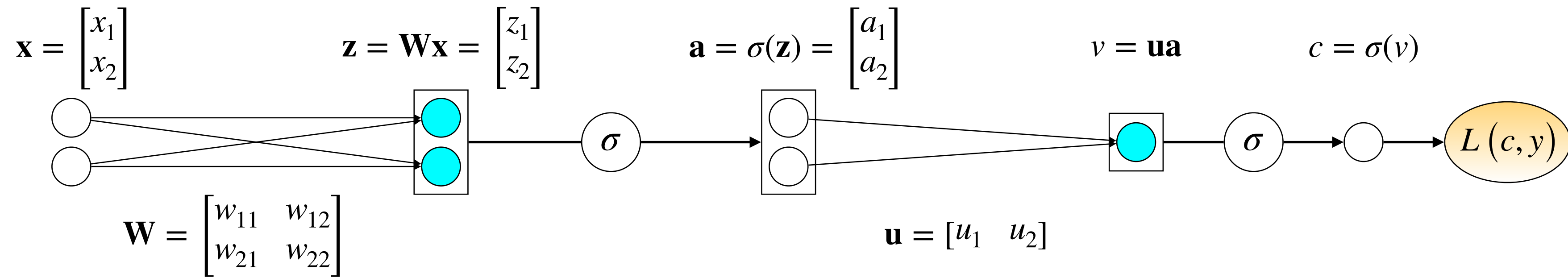


$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

$$\frac{\partial v}{\partial a_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial a_i} = u_i$$

$$\begin{aligned} \frac{\partial a_i}{\partial z_i} &= \sigma(z_i) \cdot (1 - \sigma(z_i)) \\ &= a_i \cdot (1 - a_i) \end{aligned}$$

# Backprop (toy) example (9)



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

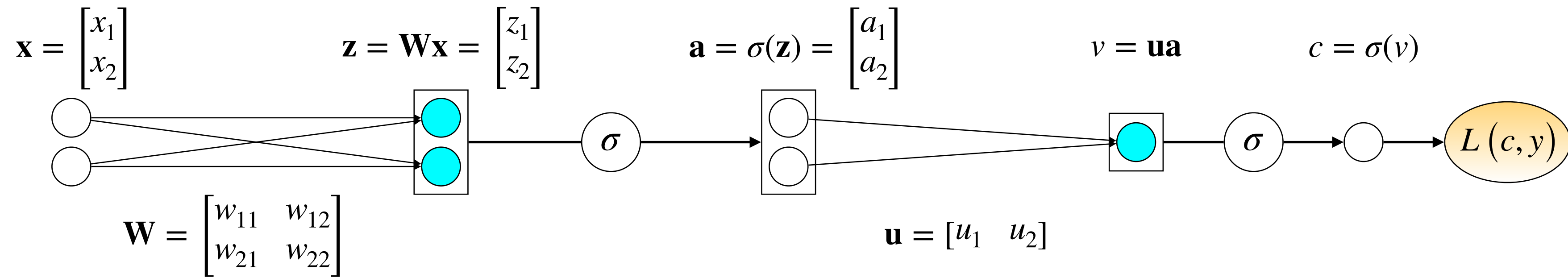
Given that  $z_i = \sum_{j=1}^2 w_{ij} \cdot x_j$

$$\frac{\partial v}{\partial a_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial a_i} = u_i$$

$$\begin{aligned} \frac{\partial a_i}{\partial z_i} &= \sigma(z_i) \cdot (1 - \sigma(z_i)) \\ &= a_i \cdot (1 - a_i) \end{aligned}$$



# Backprop (toy) example (9)



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}}$$

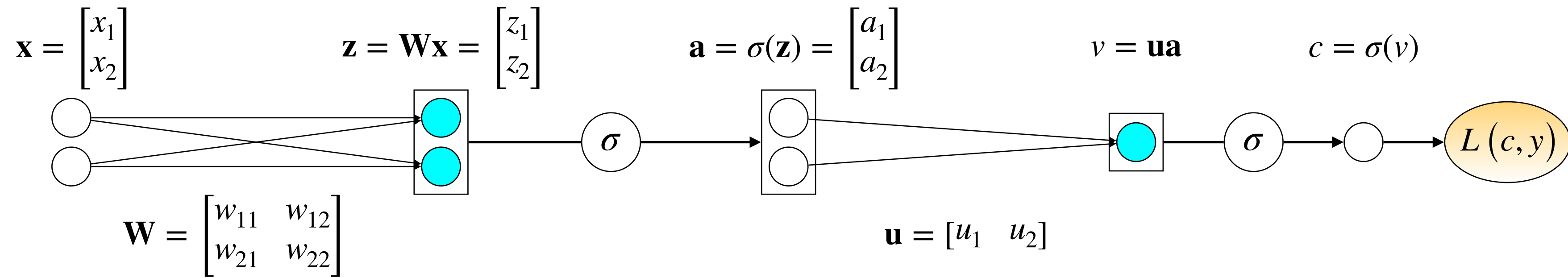
Given that  $z_i = \sum_{j=1}^2 w_{ij} \cdot x_j$

$$\frac{\partial v}{\partial a_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial a_i} = u_i$$

$$\begin{aligned} \frac{\partial a_i}{\partial z_i} &= \sigma(z_i) \cdot (1 - \sigma(z_i)) \\ &= a_i \cdot (1 - a_i) \end{aligned}$$

$$\frac{\partial z_i}{\partial w_{ij}} = x_j$$

# Backprop (toy) example (9)



$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial c} \cdot \frac{\partial c}{\partial v} \cdot \frac{\partial v}{\partial a_i} \cdot \frac{\partial a_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}} = (c - y) \cdot u_i \cdot a_i \cdot (1 - a_i) \cdot x_j$$

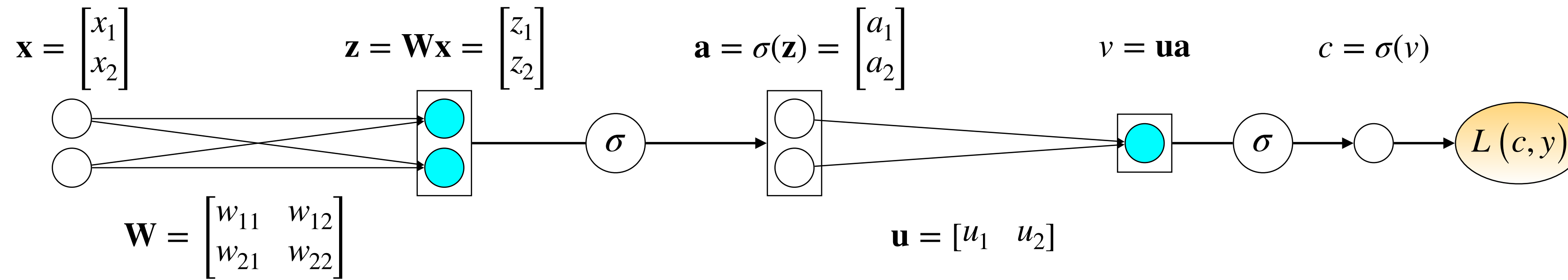
$$\frac{\partial v}{\partial a_i} = \frac{\partial \left( \sum_{i=1}^2 u_i \cdot a_i \right)}{\partial a_i} = u_i$$

$$\begin{aligned} \frac{\partial a_i}{\partial z_i} &= \sigma(z_i) \cdot (1 - \sigma(z_i)) \\ &= a_i \cdot (1 - a_i) \end{aligned}$$

Given that  $z_i = \sum_{j=1}^2 w_{ij} \cdot x_j$

$$\frac{\partial z_i}{\partial w_{ij}} = x_j$$

# Updating the parameters of the NN



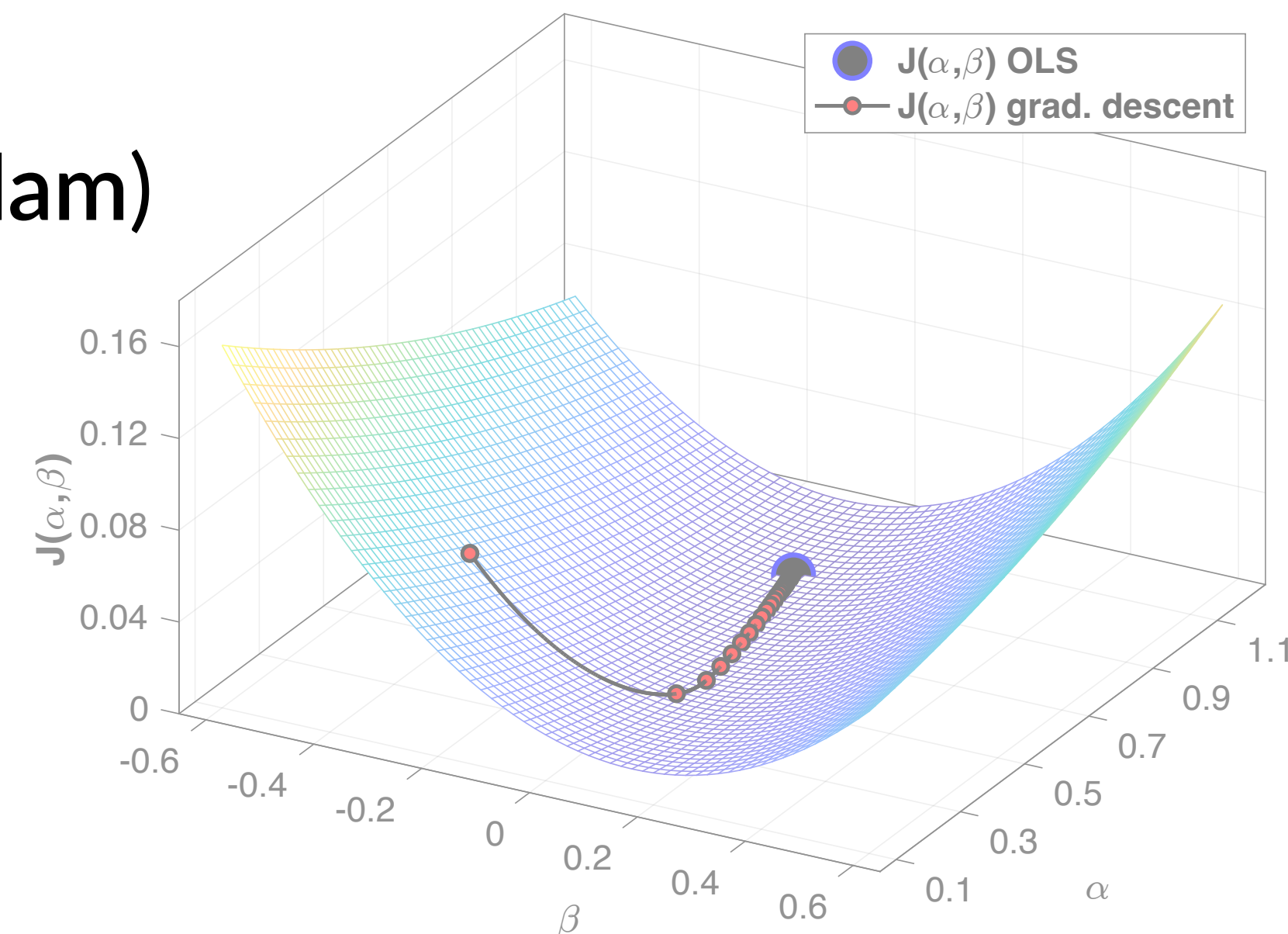
$$u_i^{\text{new}} = u_i^{\text{old}} - \eta \frac{\partial L}{\partial u_i}$$

using a learning rate  $\eta$

$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} - \eta \frac{\partial L}{\partial w_{ij}}$$

# Optimisation (training)

- ▶ Stochastic gradient descent (SGD) *works* most of the time with some effort
  - ➔ if we know our data / task well and can handle the learning rate ( $\eta$ )
- ▶ Adaptive (more sophisticated) optimisers perform generally better; *keep track how much gradients change and dynamically decide how much to update the weights*
  - ➔ RMSProp
  - ➔ Adaptive Moment Estimation Method (Adam)
  - ➔ Adagrad
  - ➔ AdaDelta
  - ➔ SparseAdam
  - ➔ many other variants



# Learning rate ( $\eta$ )

- ▶ We want the learning rate to be just right (not too large or small)
- ▶ Too large  $\implies$  learning too fast: the model may diverge and not converge
- ▶ Too small  $\implies$  learning too slow: the model will not diverge, but may take ages to converge
- ▶  $\approx 0.001$  is a common starting point value for a learning rate  
tune it by orders of magnitude e.g.  $[0.01, 0.001, 0.0001]$
- ▶ In SGD, you might want to decrease the learning rate as the training epochs increase
- ▶ In fancier optimisers (e.g. Adam) we set the initial learning rate, but then the optimiser takes care of dynamically tuning it

# Next lecture

- ▶ Friday, January 26
- ▶ Word embeddings (`word2vec` mainly)

