

Neural Networks

COMP3314 — Lecture 5

Lingpeng Kong

Department of Computer Science, The University of Hong Kong

Based on: Probabilistic Machine Learning by Kevin Murphy

Slides from: Saw Shier Nee with special thanks!

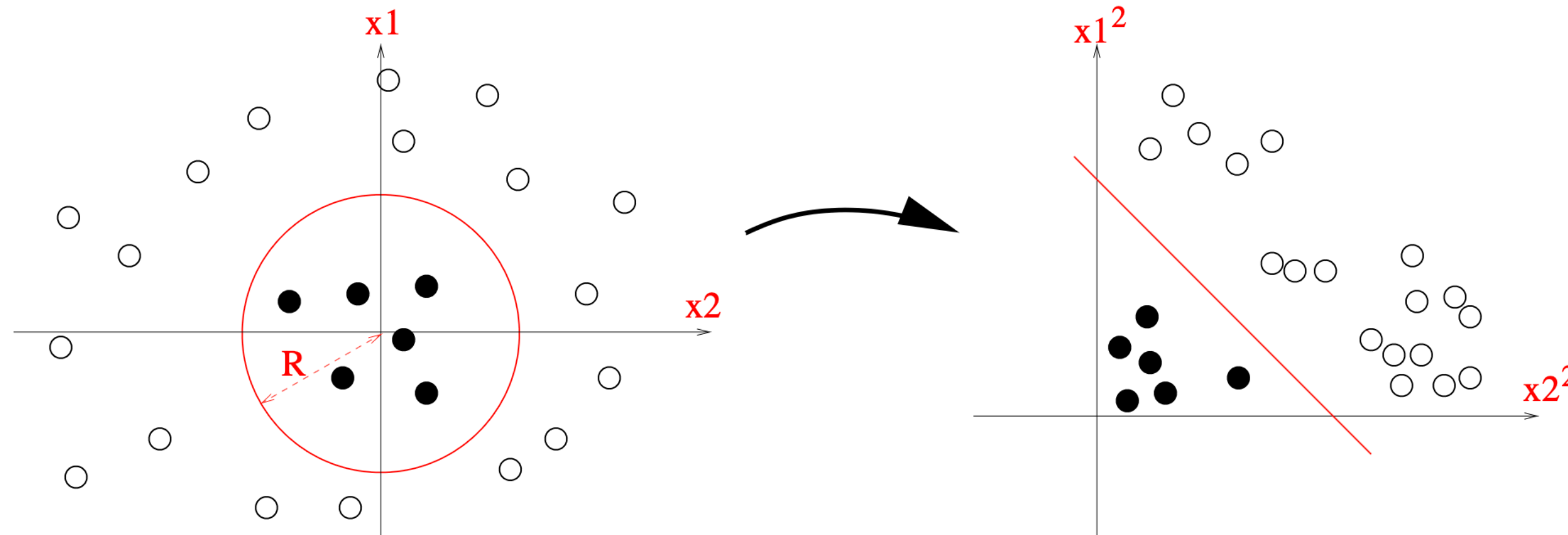
Basis Function Expansion

$$p(y|\mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y|\sigma(\mathbf{w}^\top \mathbf{x} + b))$$

Logistic Regression

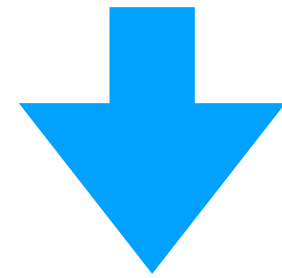
$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\boldsymbol{\phi}(\mathbf{x}) + b$$

$$\boldsymbol{\phi}(\mathbf{x}) = [1, x, x^2, x^3, \dots]$$

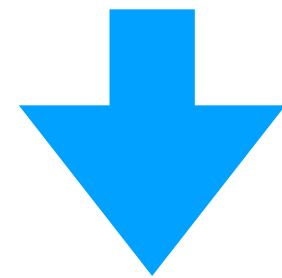


Basis Function Expansion

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\phi(\mathbf{x}) + \mathbf{b} \quad \phi(x) = [1, x, x^2, x^3, \dots]$$



$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\phi(\mathbf{x}; \boldsymbol{\theta}_2) + \mathbf{b}$$



$$f(\mathbf{x}; \boldsymbol{\theta}) = f_L(f_{L-1}(\dots(f_1(\mathbf{x}))\dots)) \quad f_\ell(\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}_\ell)$$

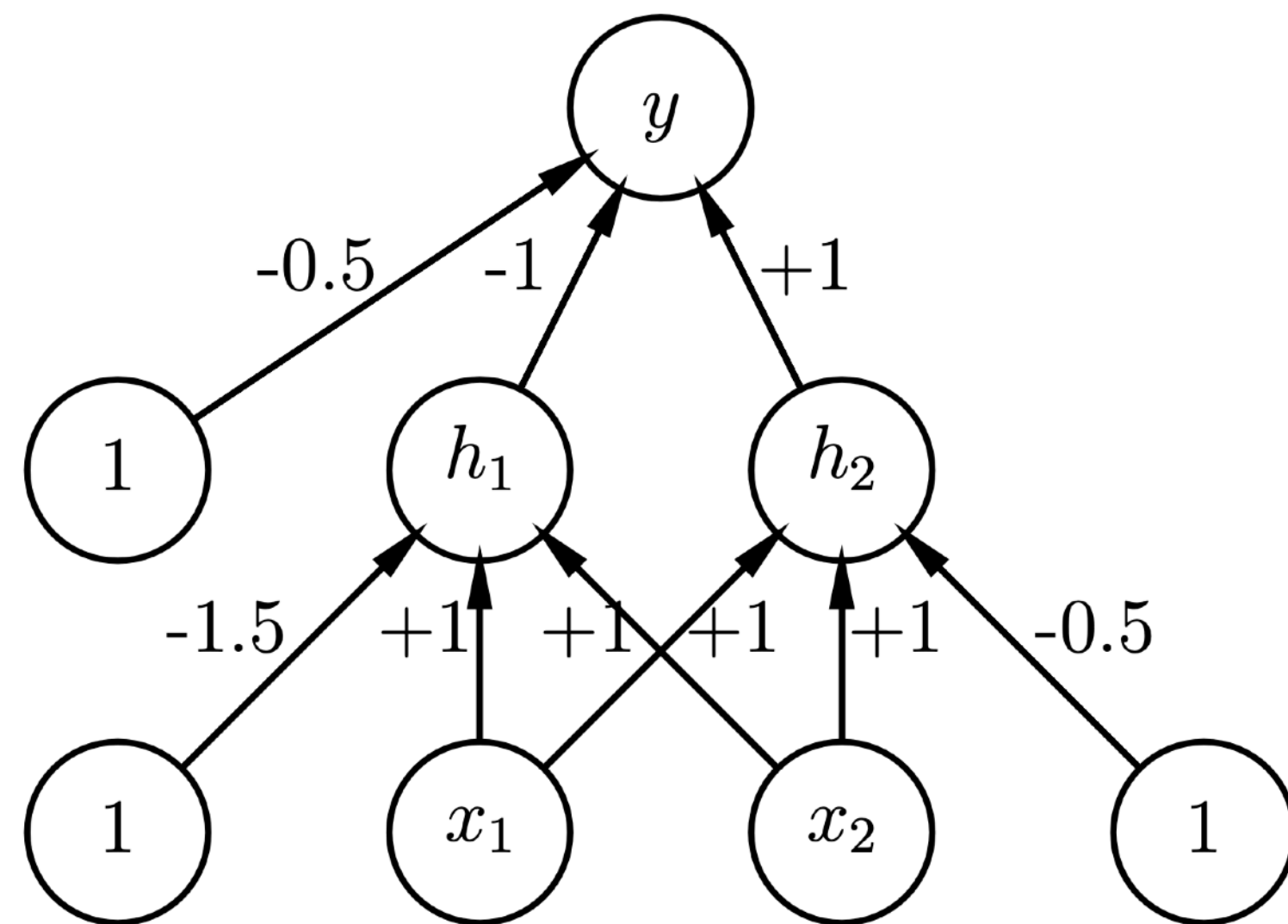
Deep Neural Networks (DNNs)

Multilayer Perceptrons (MLPs)

perceptron: a deterministic version of logistic regression

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbb{I}(\mathbf{w}^T \mathbf{x} + b \geq 0) = H(\mathbf{w}^T \mathbf{x} + b)$$

heaviside step function / linear threshold function



x1	x2	XOR(x1, x2)
0	0	0
0	1	1
1	0	1
1	1	0

Multilayer Perceptrons (MLPs)

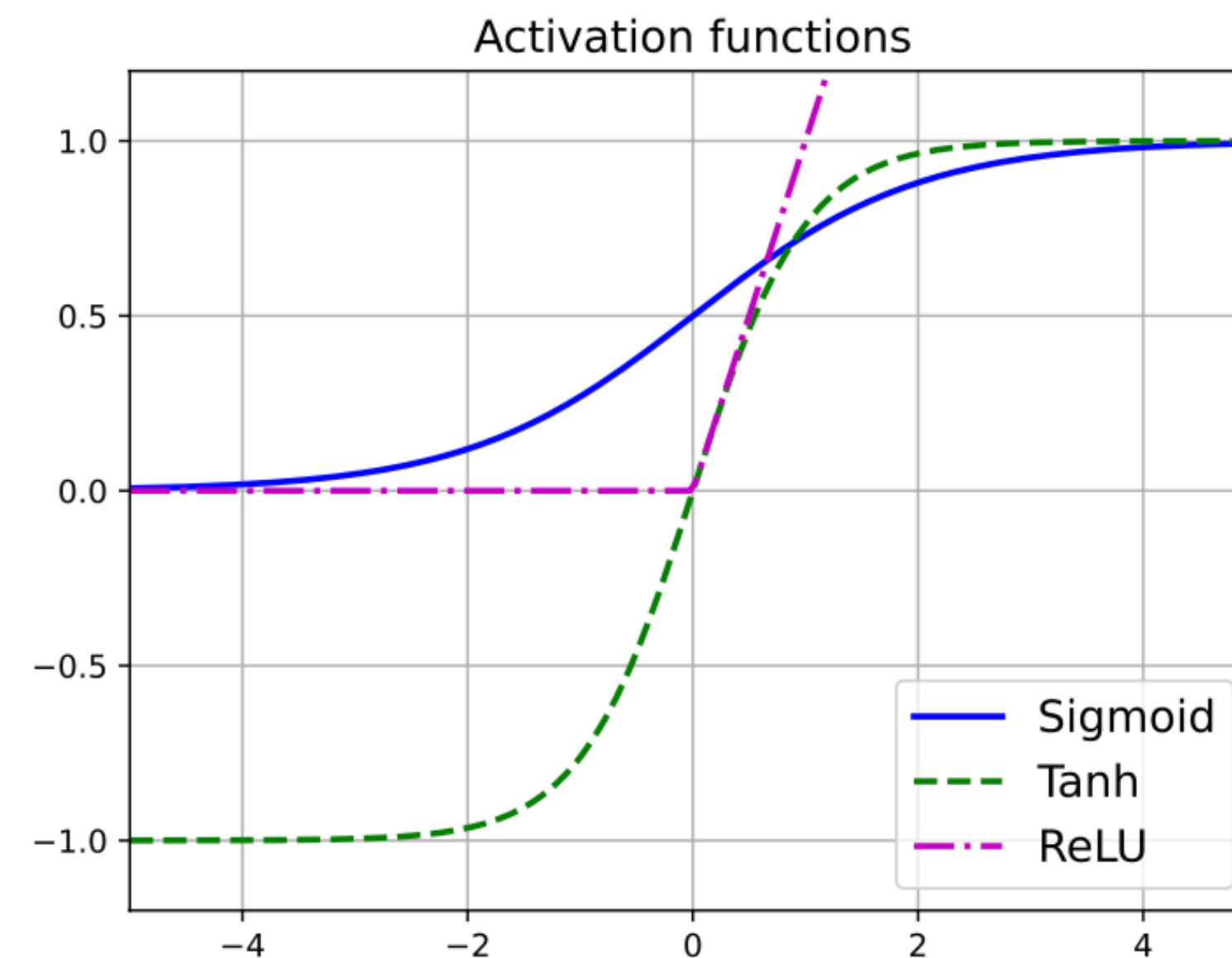
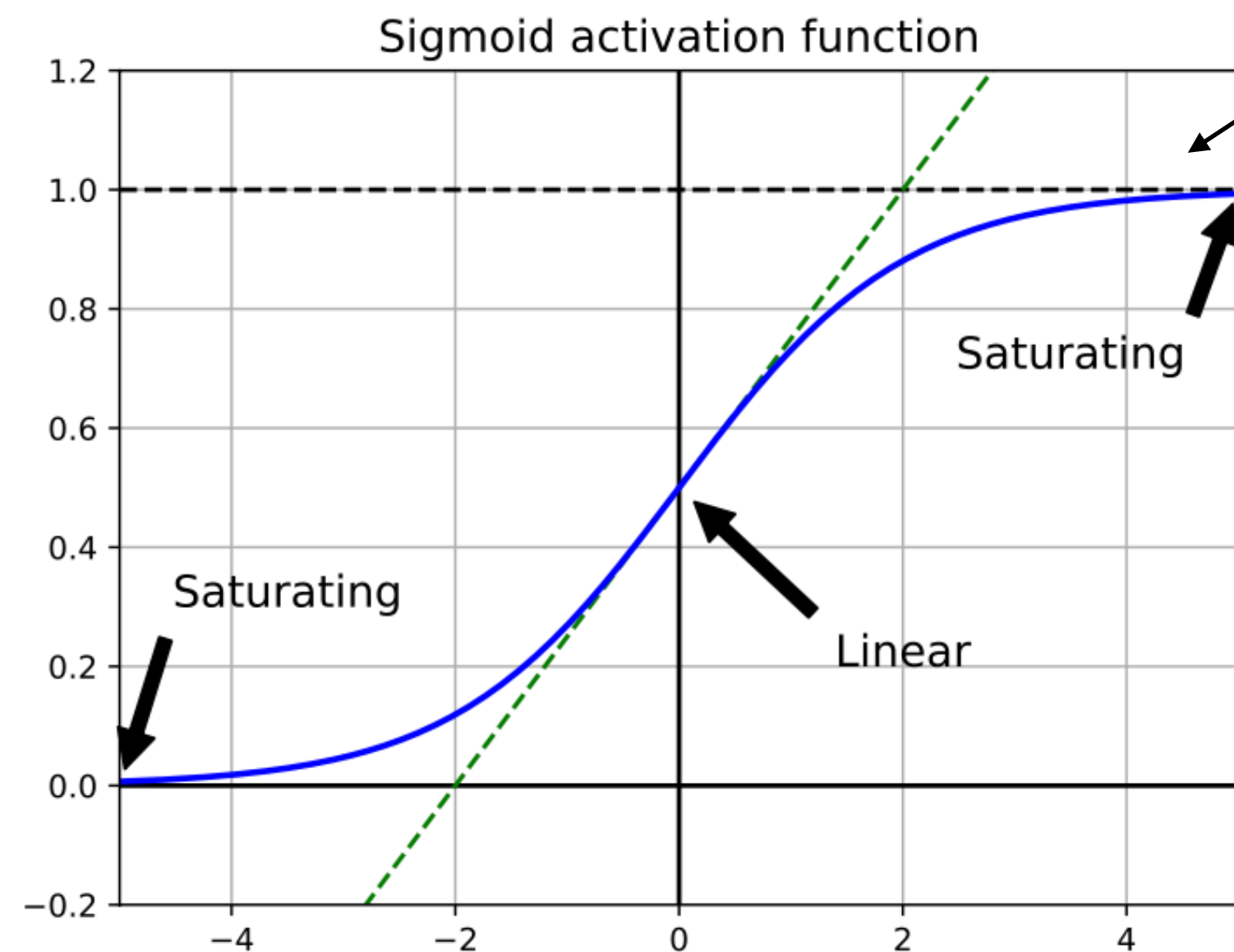
Differentiable MLPs

$$\mathbf{z}_l = f_l(\mathbf{z}_{l-1}) = \varphi_l(\mathbf{b}_l + \mathbf{W}_l \mathbf{z}_{l-1})$$

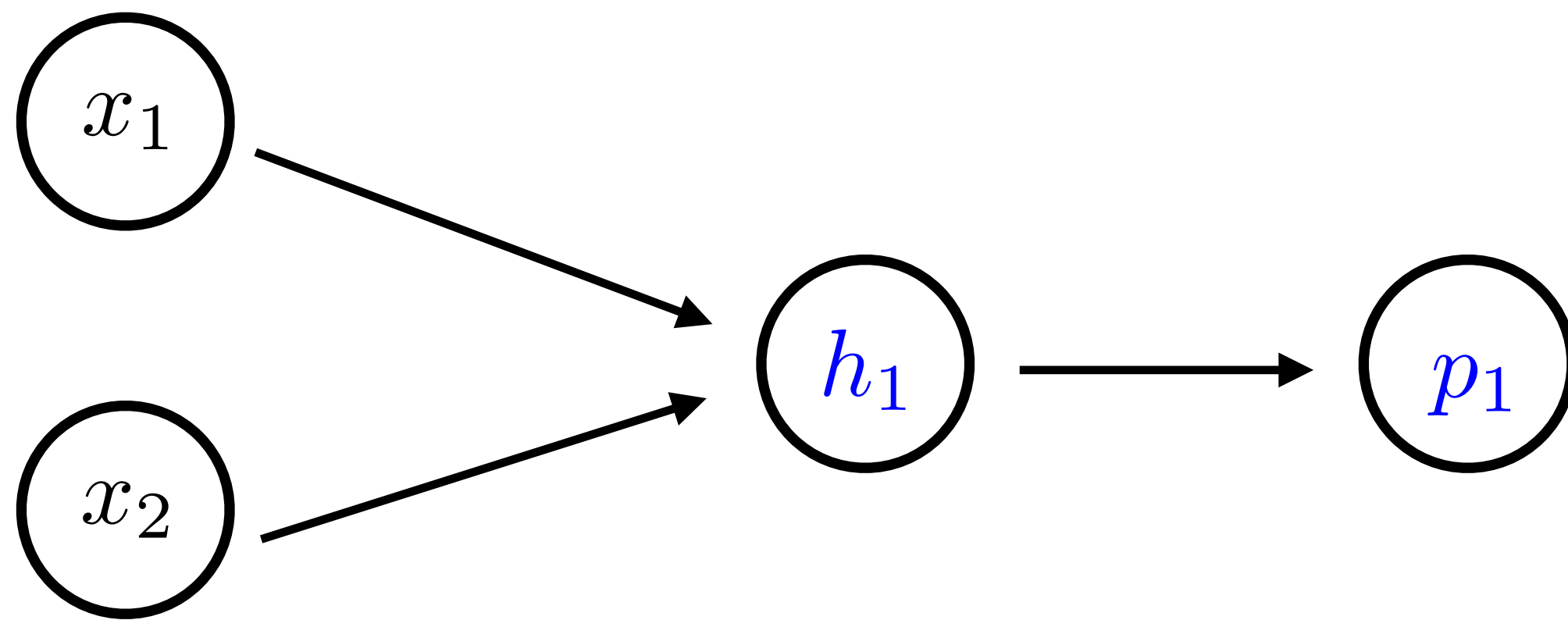
Activation function

What's bad here?

$\text{ReLU}(a) = \max(a, 0) = a\mathbb{I}(a > 0)$
(ReLU: rectified linear unit)



Computational Graphs

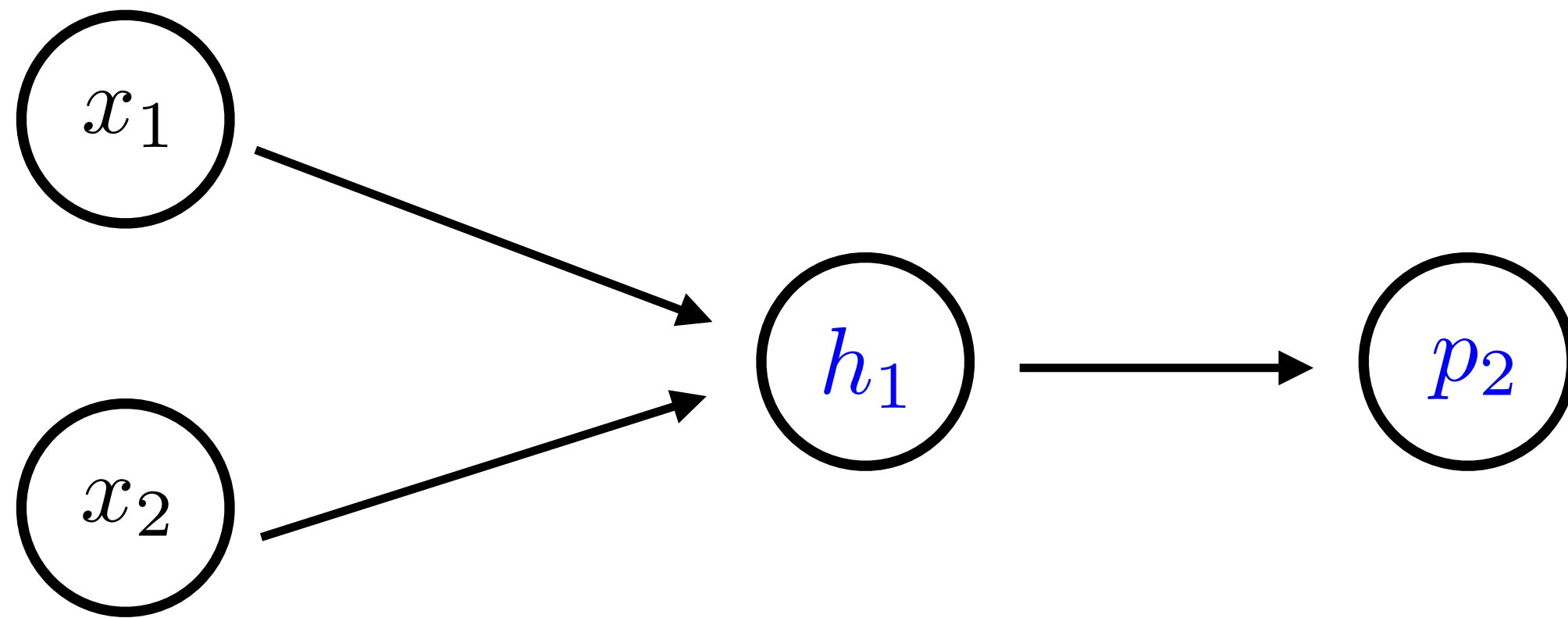


$$p_1 := p(y = 1 \mid x_1, x_2)$$

$$h_1 = w_1 x_1 + w_2 x_2 + b$$

$$p_1 = \frac{1}{1 + \exp(-h_1)}$$

Logistic Regression



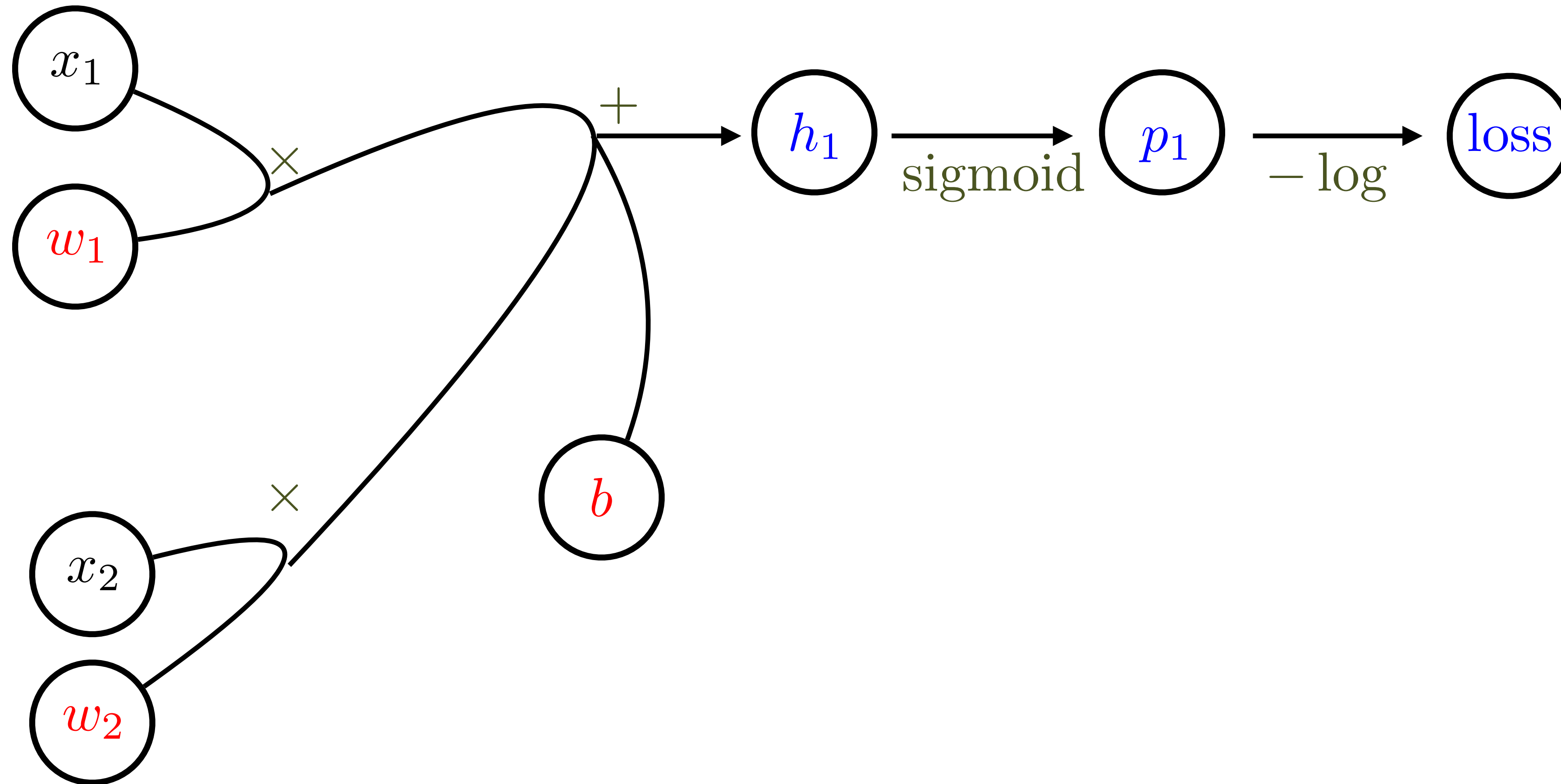
$$p_2 := p(y = 0 \mid x_1, x_2)$$

$$h_1 = w_1 x_1 + w_2 x_2 + b$$

$$p_2 = 1 - \frac{1}{1 + \exp(-h_1)} = \frac{\exp(-h_1)}{1 + \exp(-h_1)}$$

Loss Function

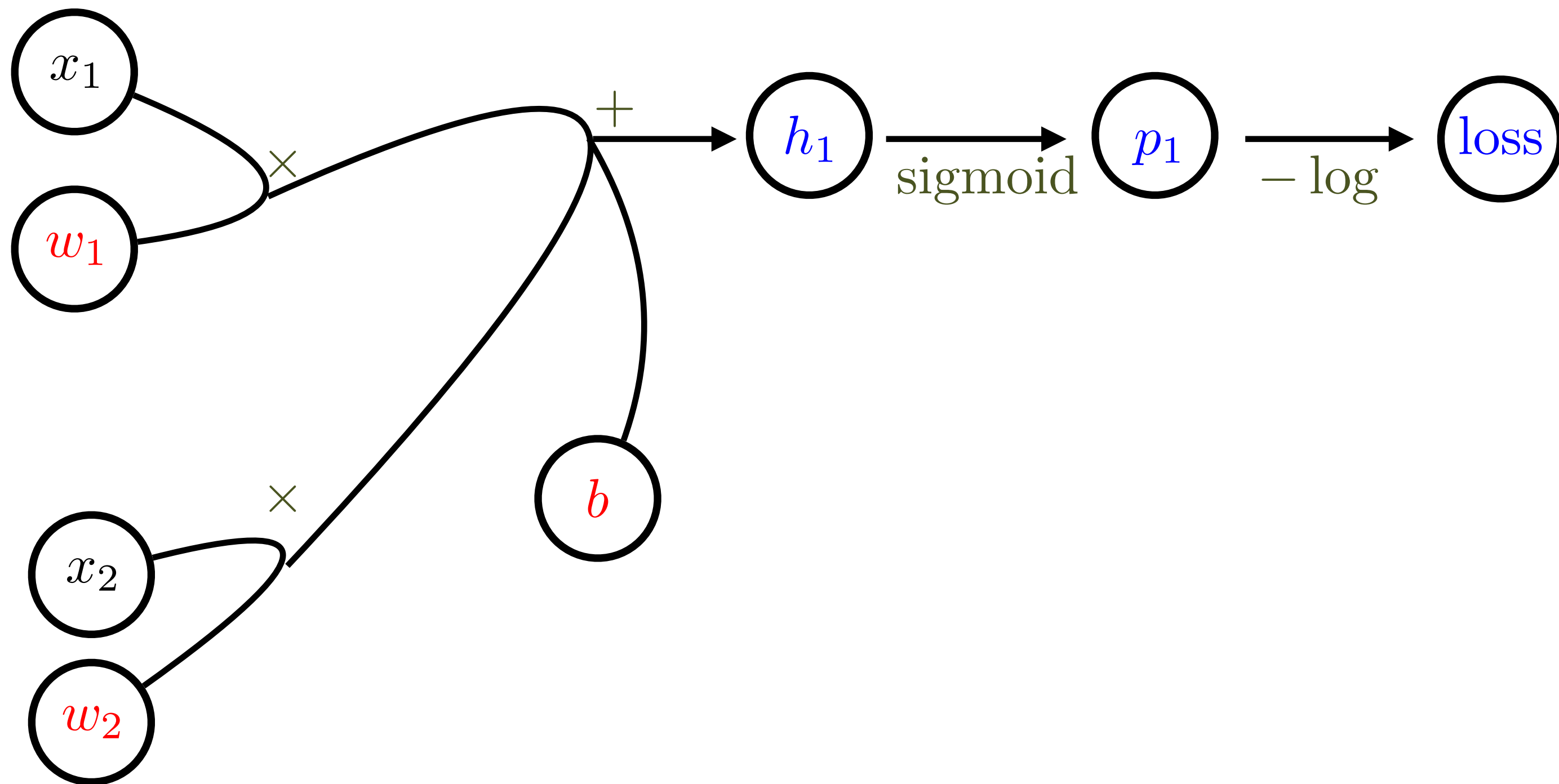
case $y = 1$:



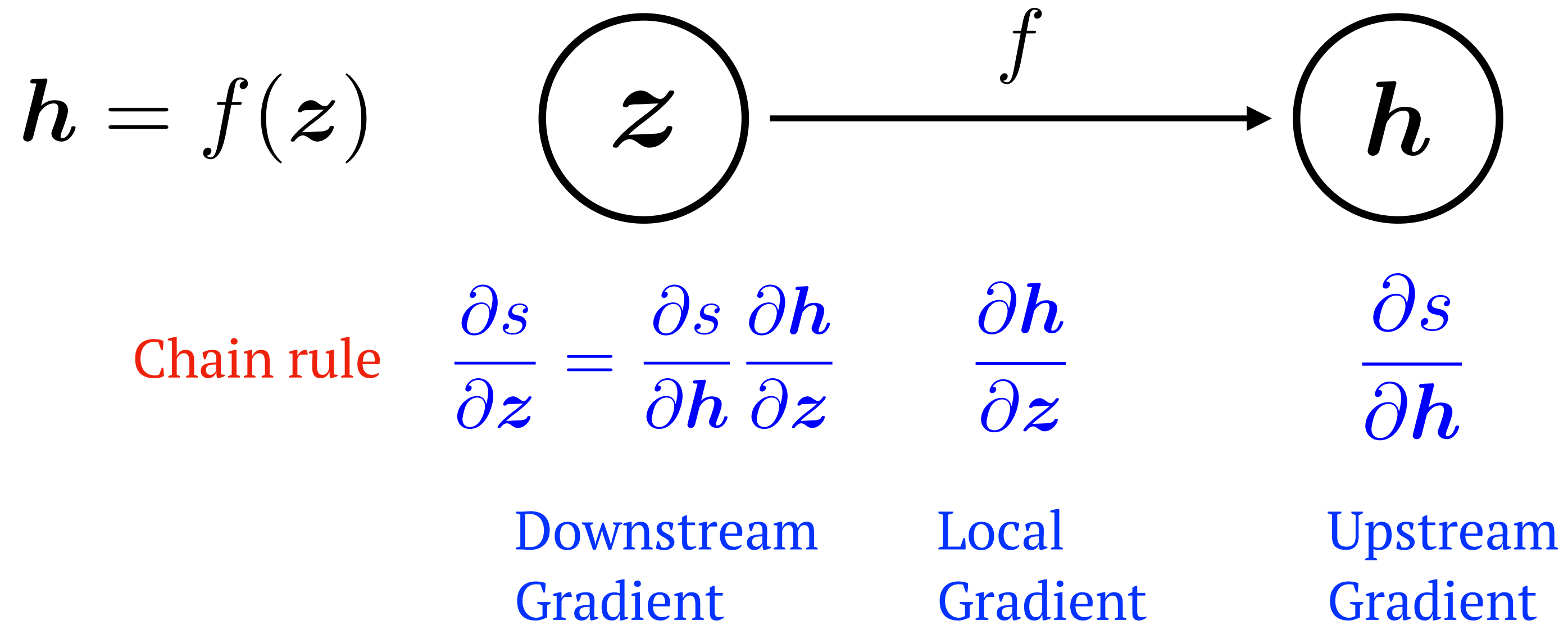
Computational Graphs

Input	x_1	x_2	
Parameter	w_1	w_2	
Expression	h_1	p_1	loss
Operation	\times	$+$	sigmoid $-\log$

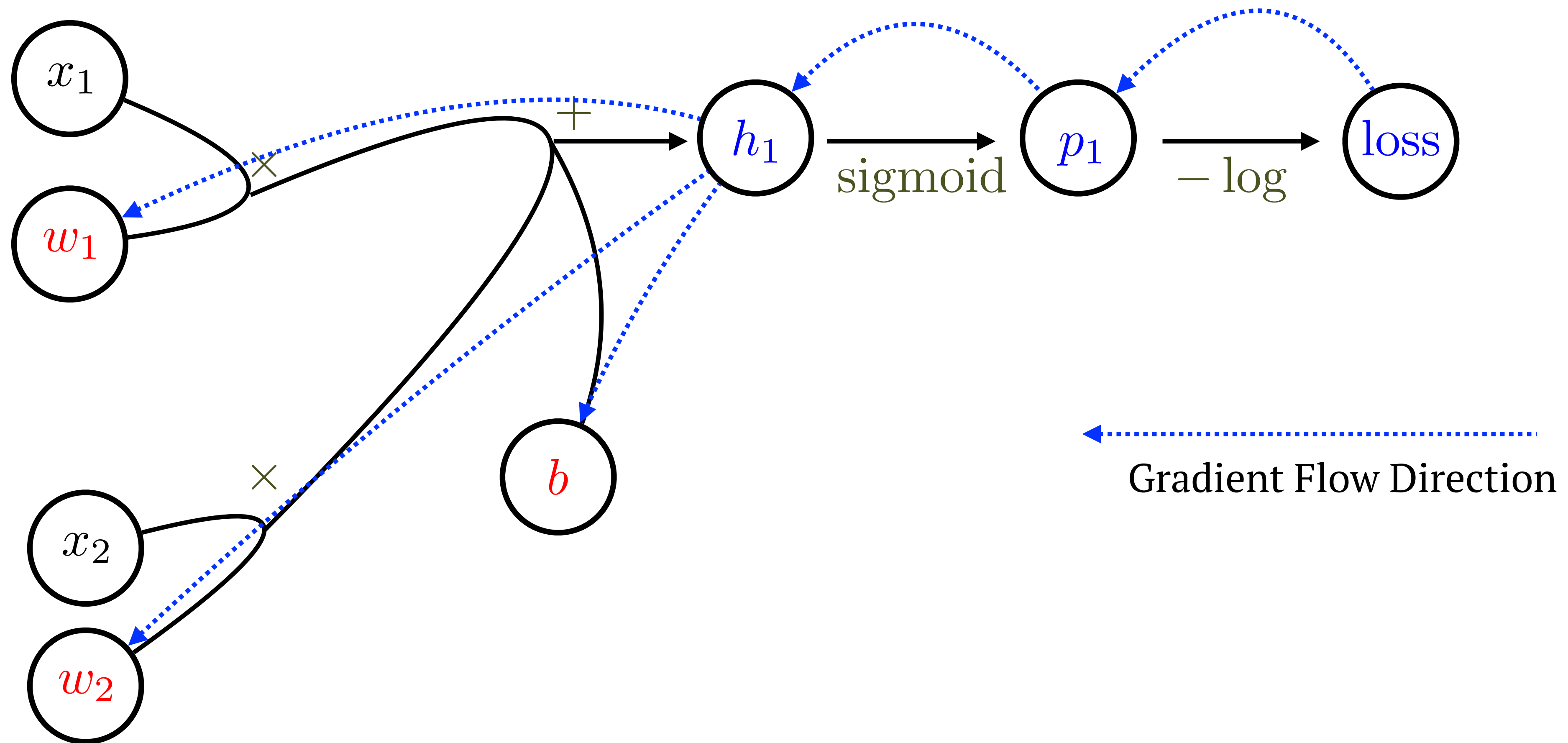
Special



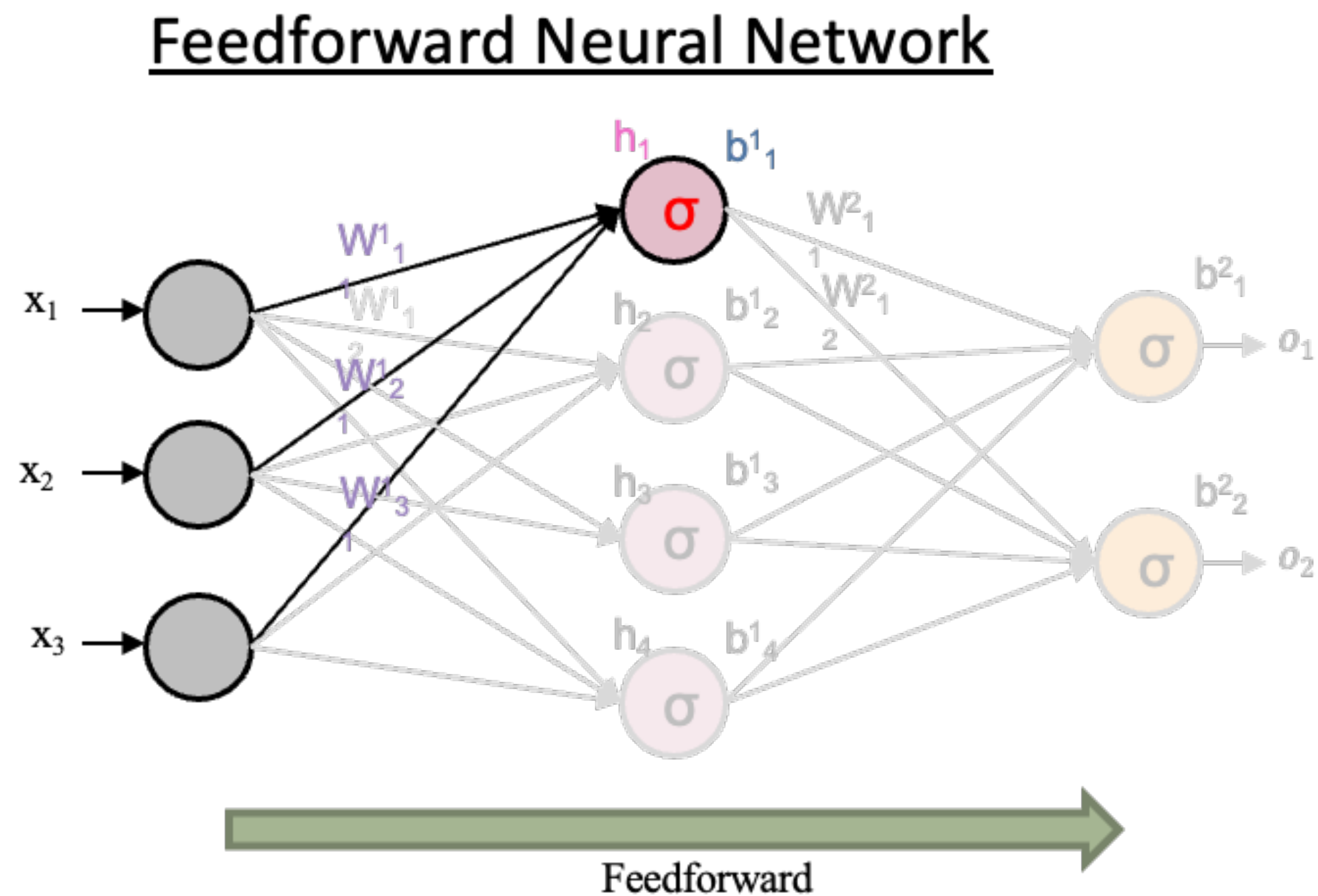
How to minimize? (Automatic Differentiation)



How to minimize? (Automatic Differentiation)



Example: Feedforward Neural Network

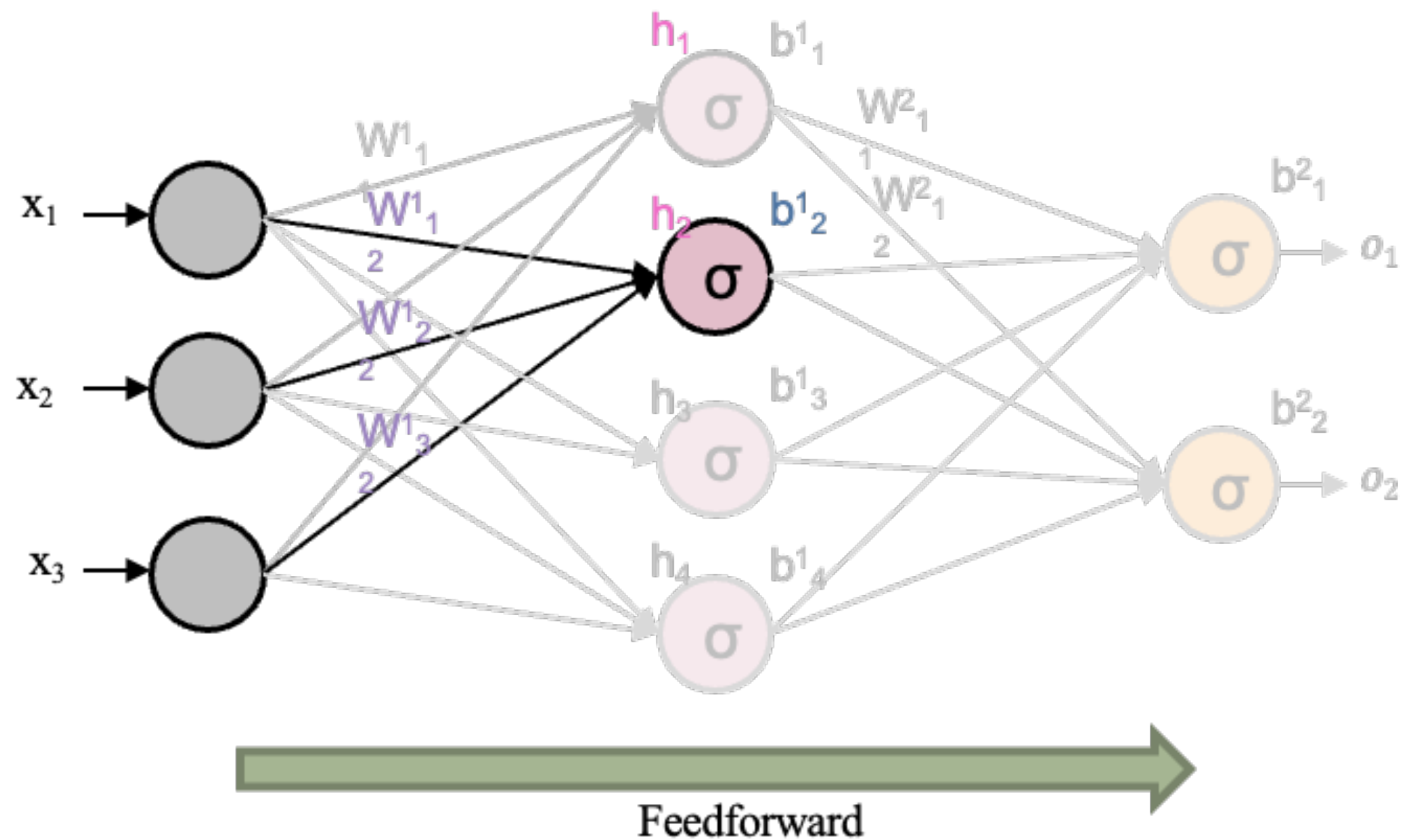


Input x Weight + Bias

$$h_1 = \sigma(x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b_1^1)$$

↑
Activation function

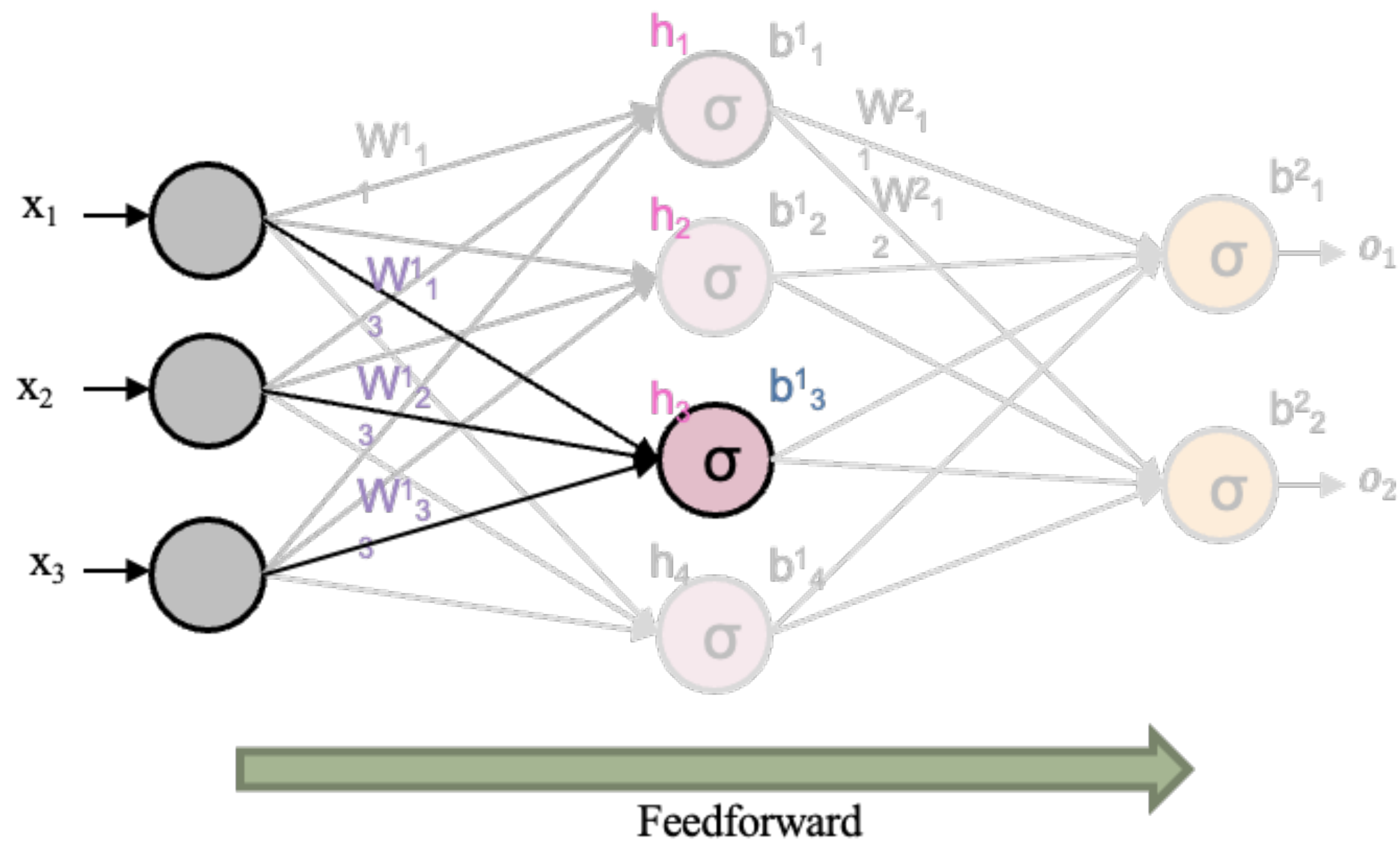
Example: Feedforward Neural Network



$$h_1 = \sigma(x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b_1^1)$$

$$h_2 = \sigma(x_1 w_{12}^1 + x_2 w_{22}^1 + x_3 w_{32}^1 + b_2^1)$$

Example: Feedforward Neural Network

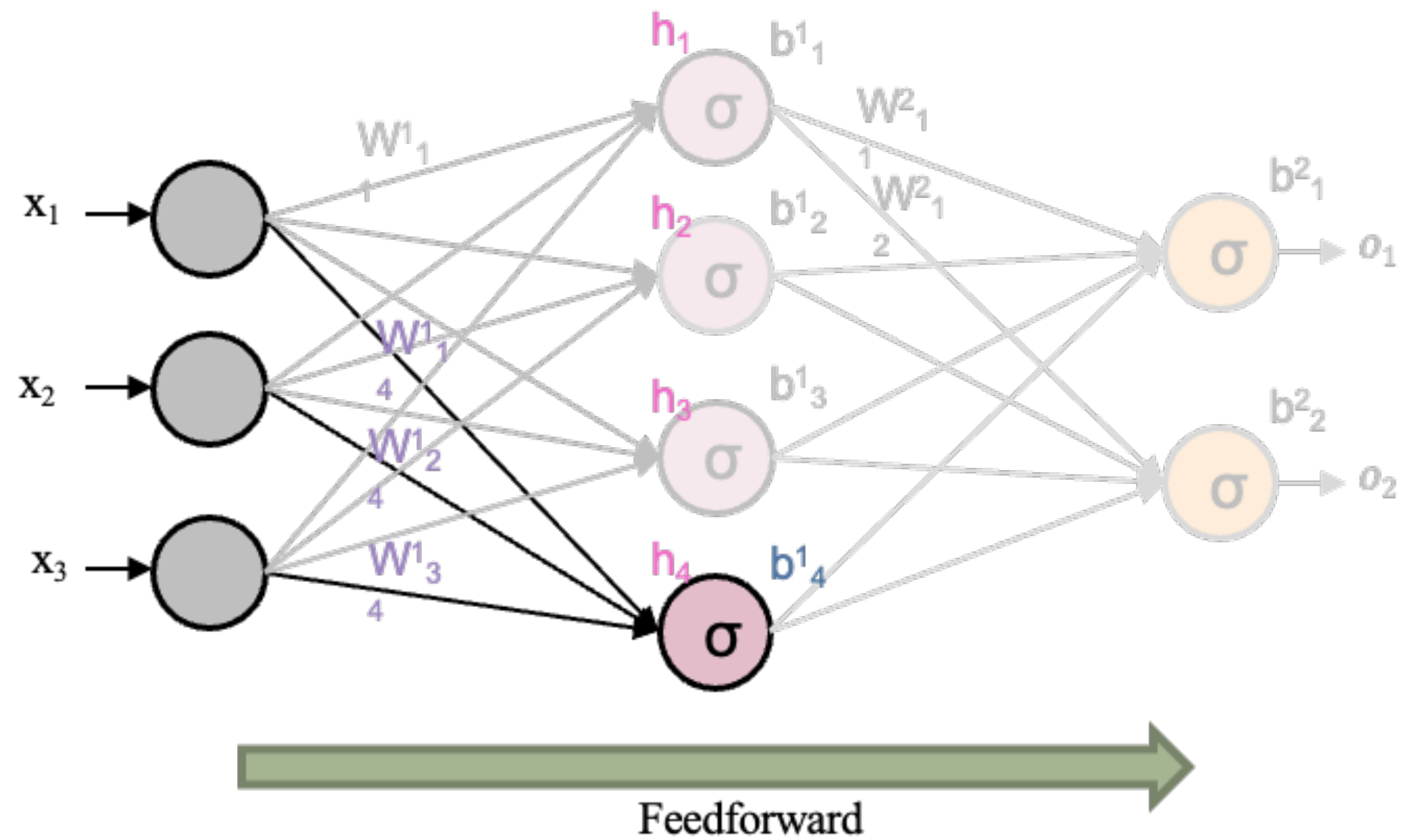


$$h_1 = \sigma(x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b_1^1)$$

$$h_2 = \sigma(x_1 w_{12}^1 + x_2 w_{22}^1 + x_3 w_{32}^1 + b_2^1)$$

$$h_3 = \sigma(x_1 w_{13}^1 + x_2 w_{23}^1 + x_3 w_{33}^1 + b_3^1)$$

Example: Feedforward Neural Network



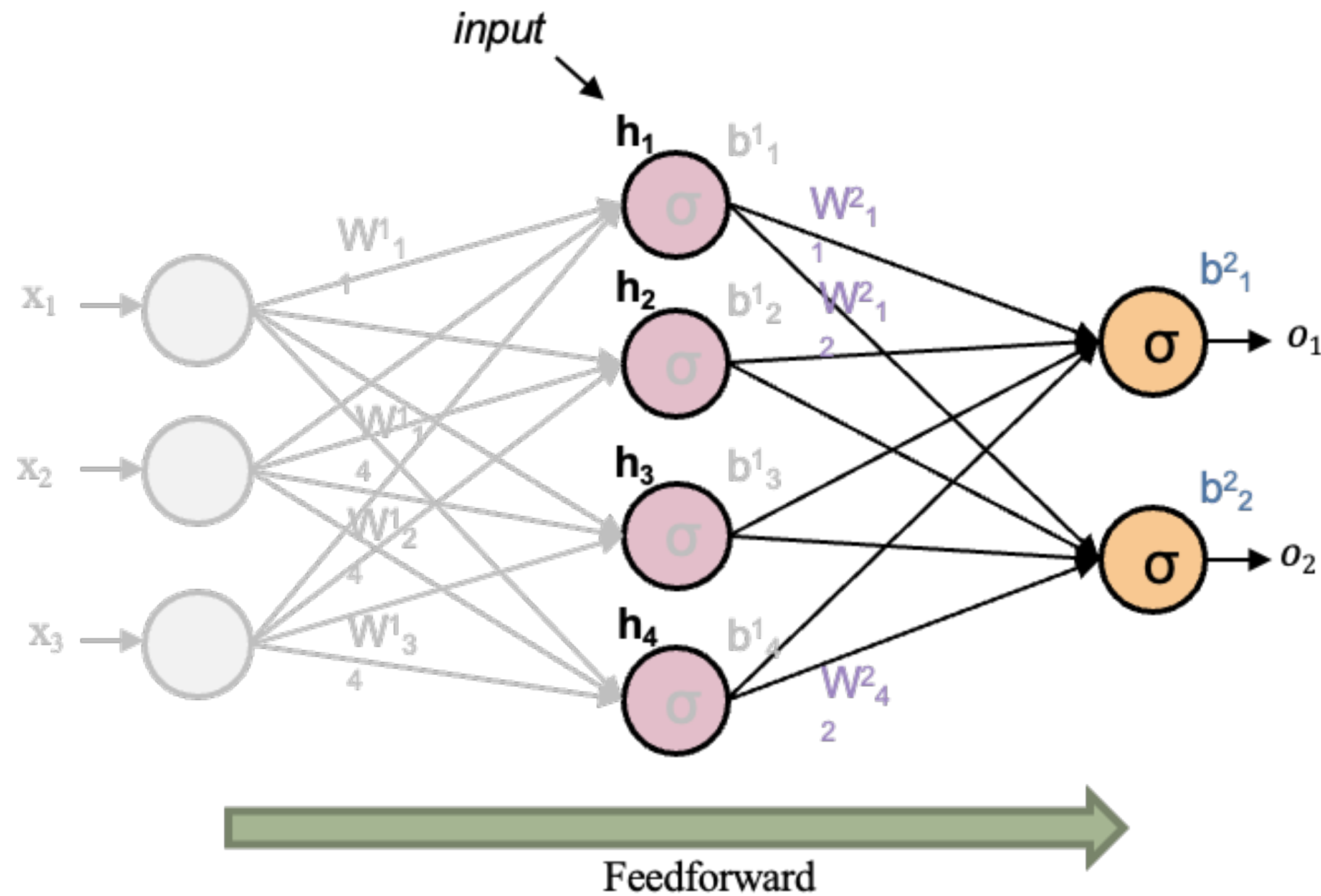
$$h_1 = \sigma(x_1 w_{11}^1 + x_2 w_{21}^1 + x_3 w_{31}^1 + b_1^1)$$

$$h_2 = \sigma(x_1 w_{12}^1 + x_2 w_{22}^1 + x_3 w_{32}^1 + b_2^1)$$

$$h_3 = \sigma(x_1 w_{13}^1 + x_2 w_{23}^1 + x_3 w_{33}^1 + b_3^1)$$

$$h_4 = \sigma(x_1 w_{14}^1 + x_2 w_{24}^1 + x_3 w_{34}^1 + b_4^1)$$

Example: Feedforward Neural Network

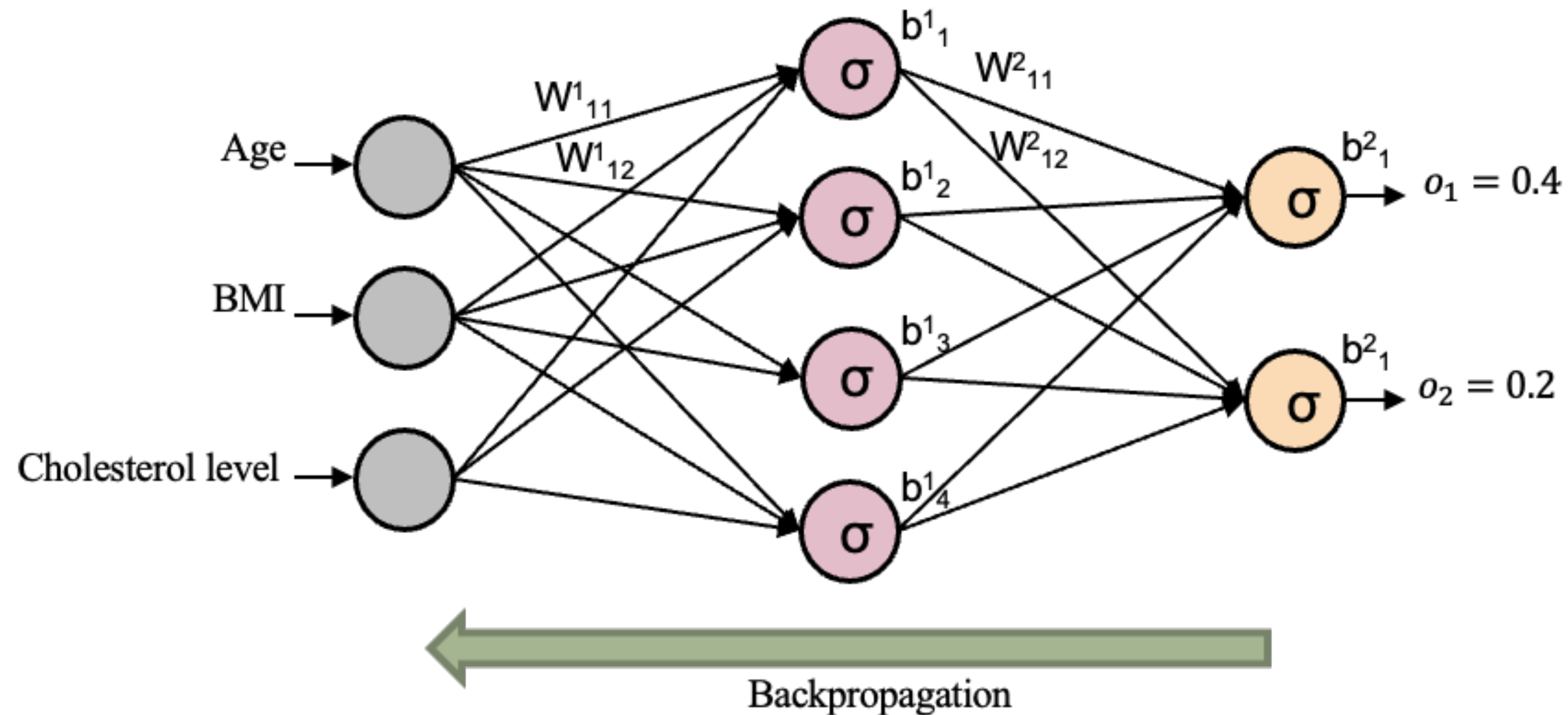


$$o_1 = \sigma(h_1 w_{11}^2 + h_2 w_{21}^2 + h_3 w_{31}^2 + h_4 w_{41}^2 + b_1^2)$$

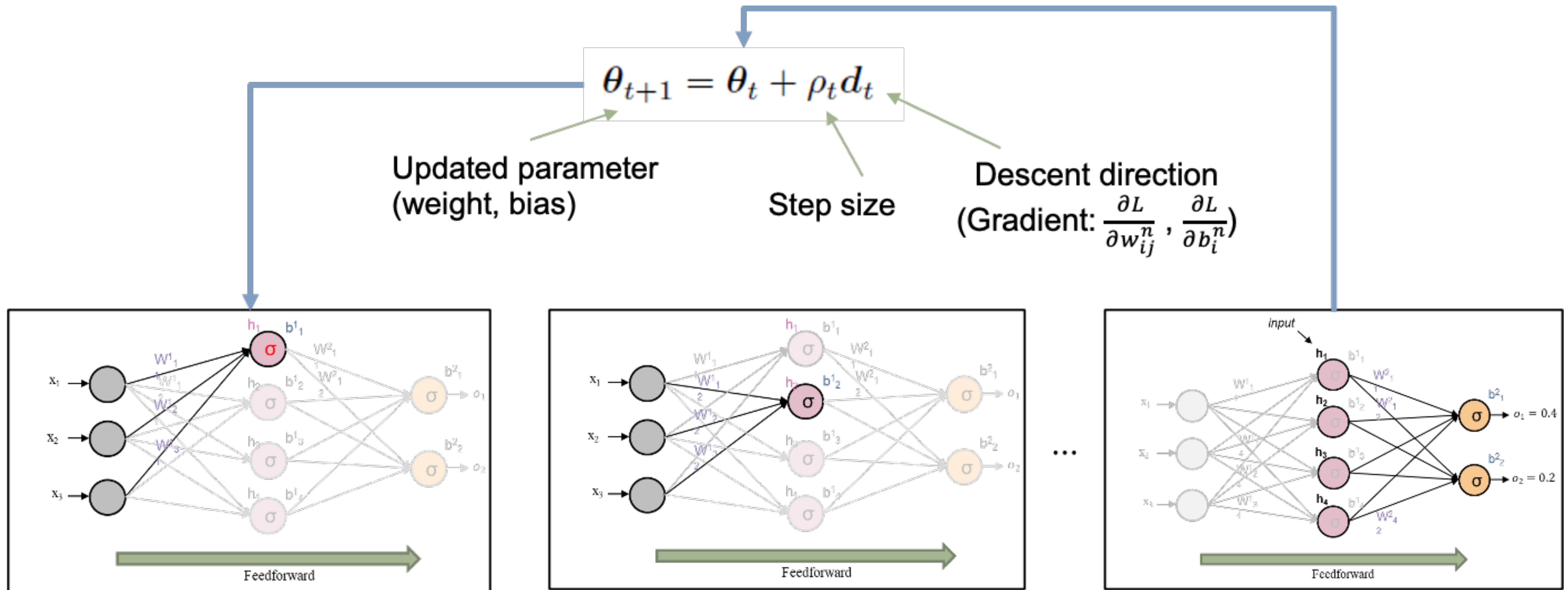
$$o_2 = \sigma(h_1 w_{12}^2 + h_2 w_{22}^2 + h_3 w_{32}^2 + h_4 w_{42}^2 + b_2^2)$$

Example: Feedforward Neural Network

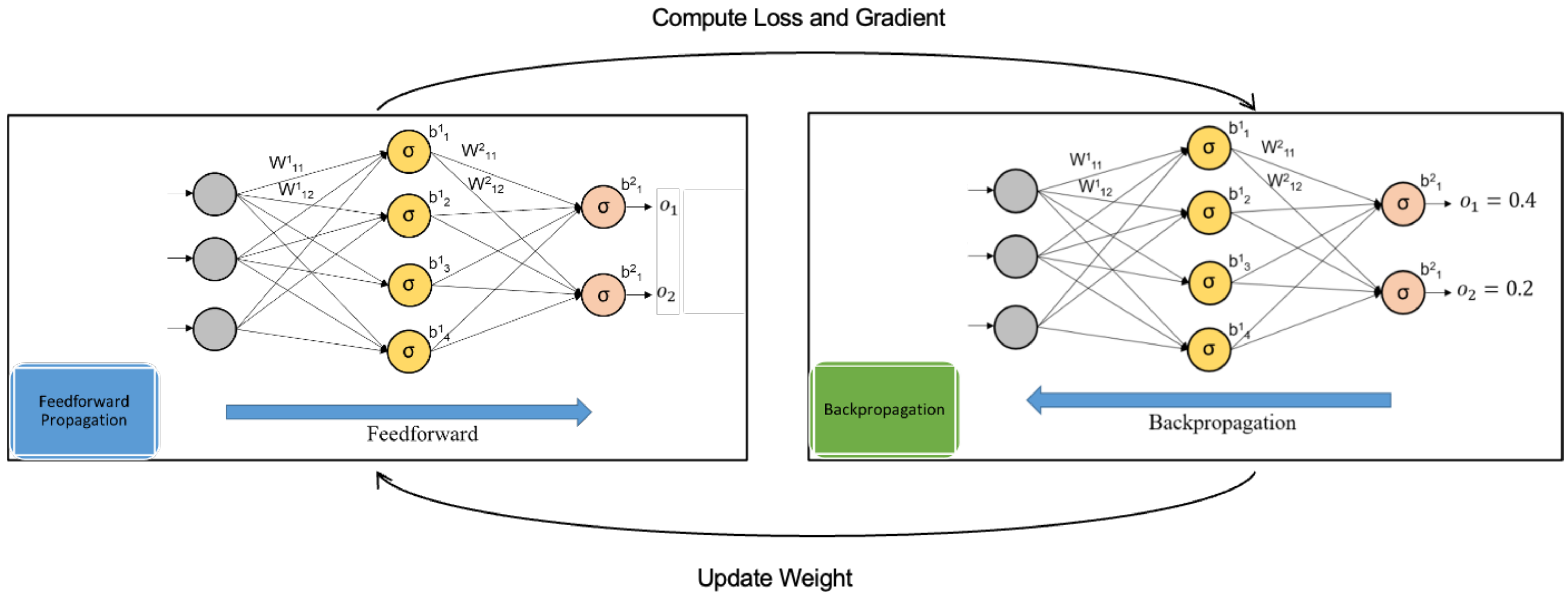
- Compute Gradients, i.e: $\frac{\partial L}{\partial w_{ij}^n}$, $\frac{\partial L}{\partial b_i^n}$
- Backpropagate the gradient to updates the weights



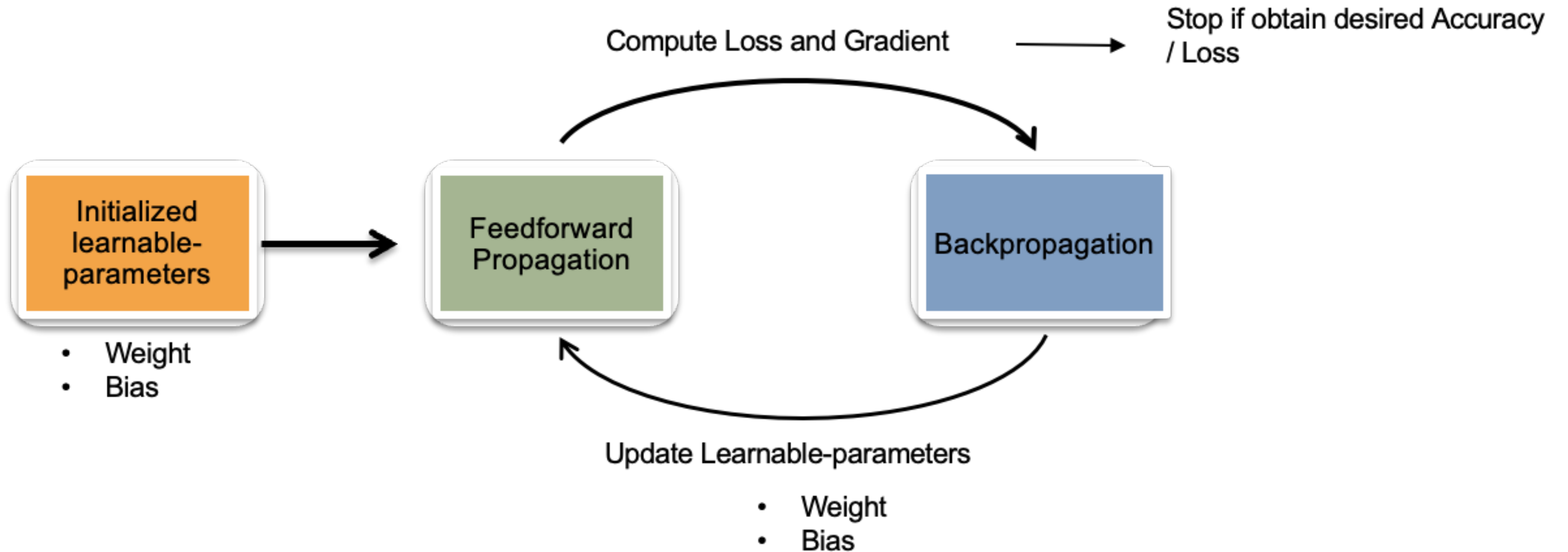
Example: Feedforward Neural Network



Example: Feedforward Neural Network



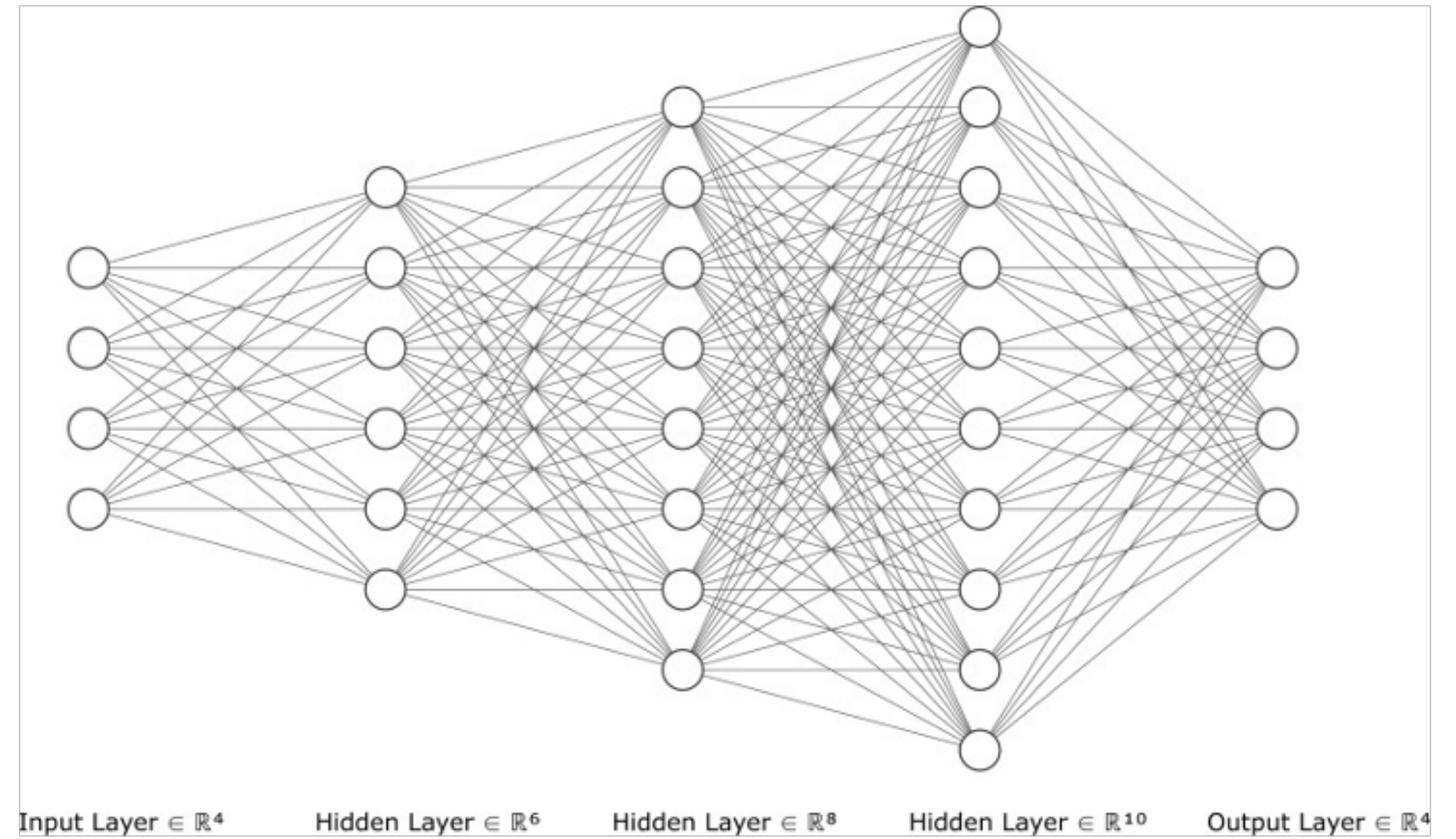
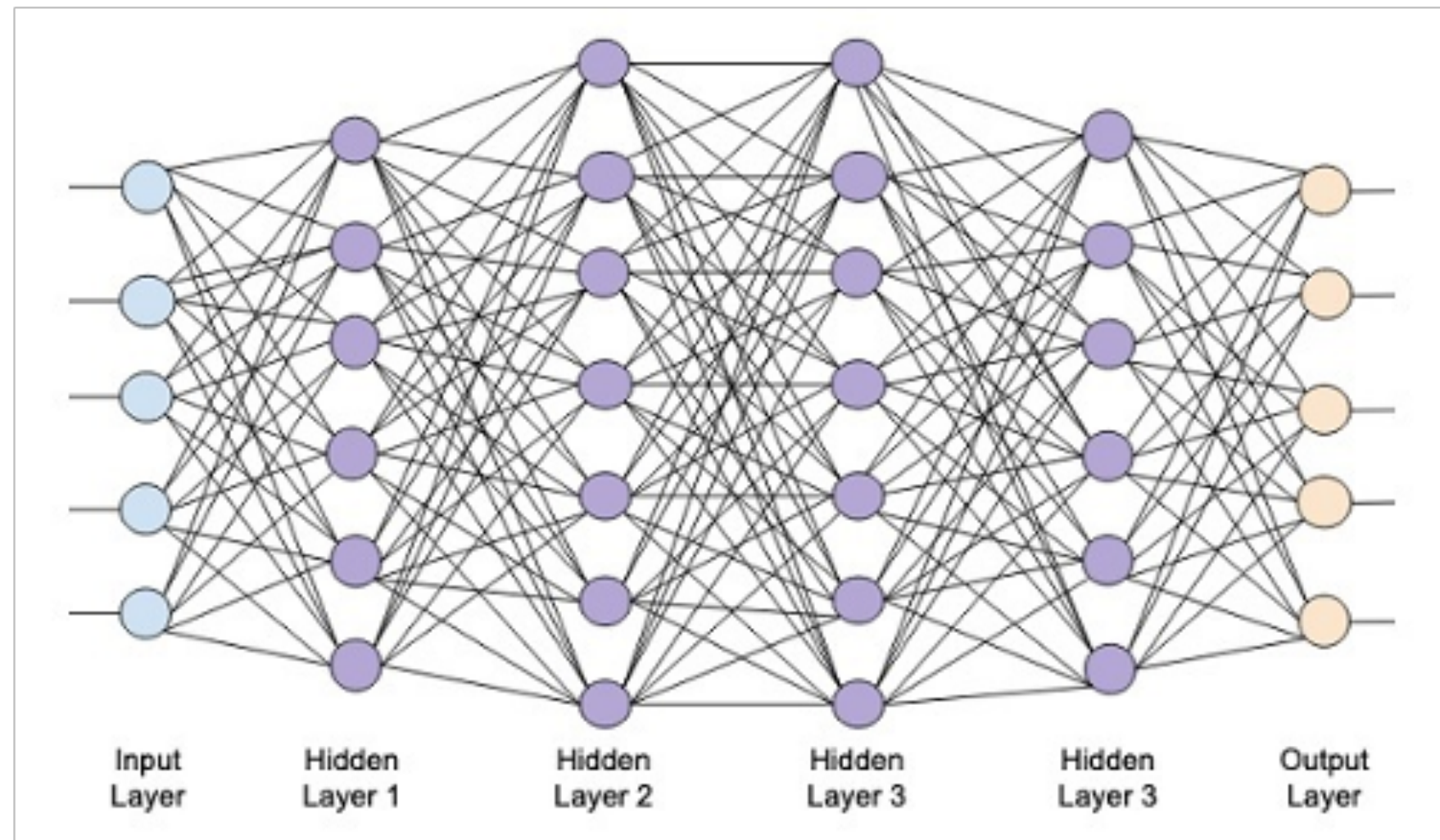
Summary



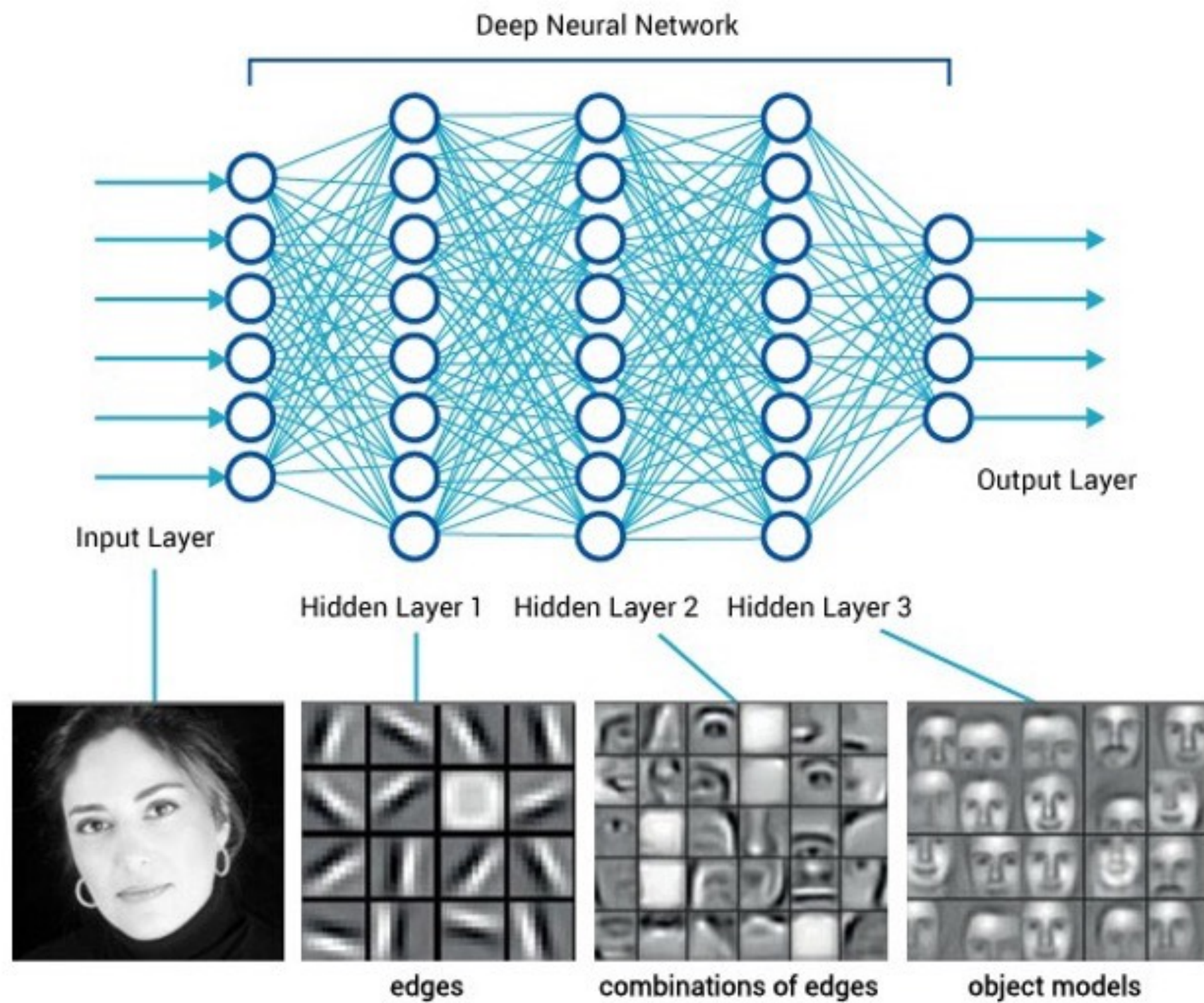
Larger Networks

Increase number of nodes

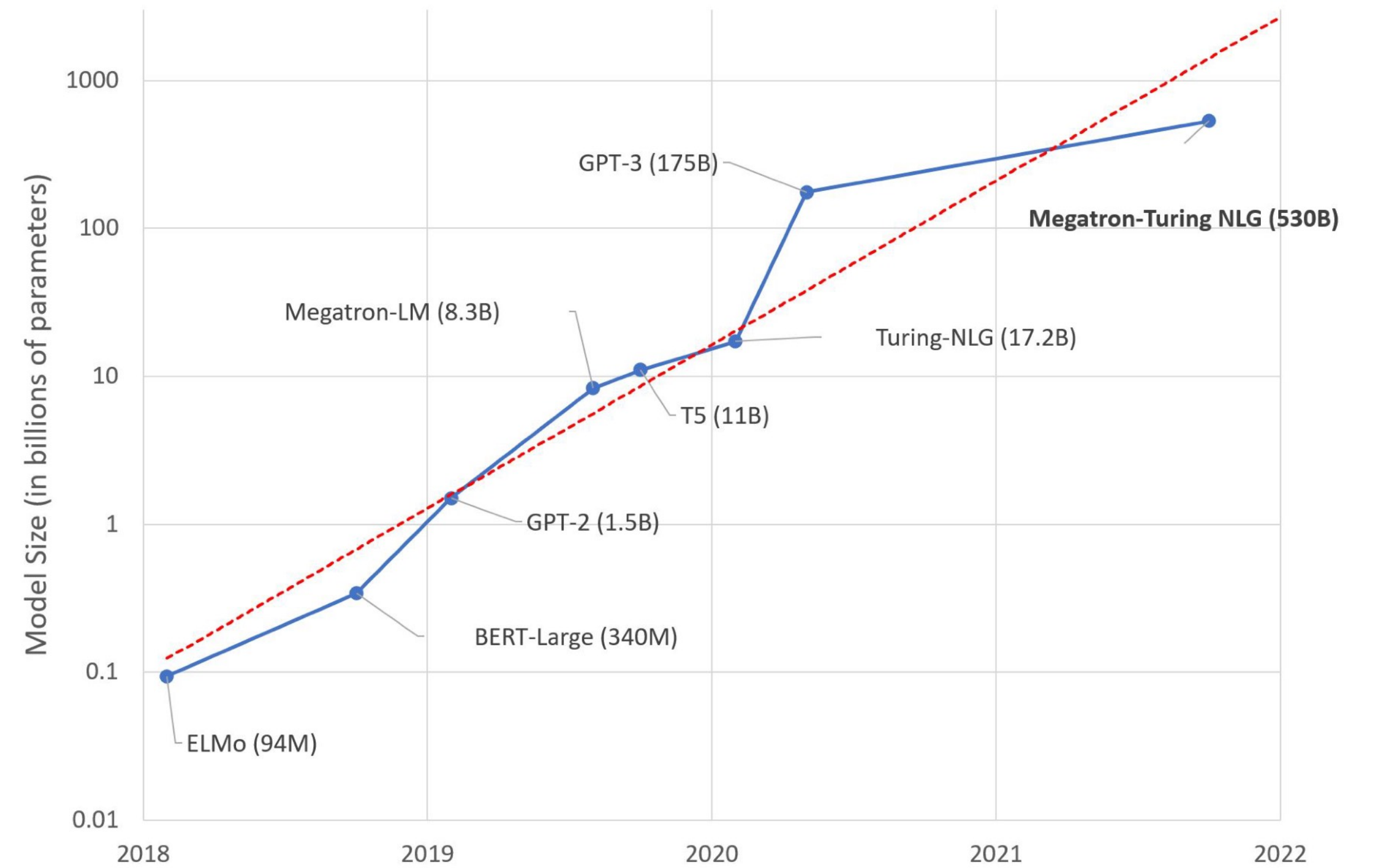
Increase number of hidden layers



Larger Networks



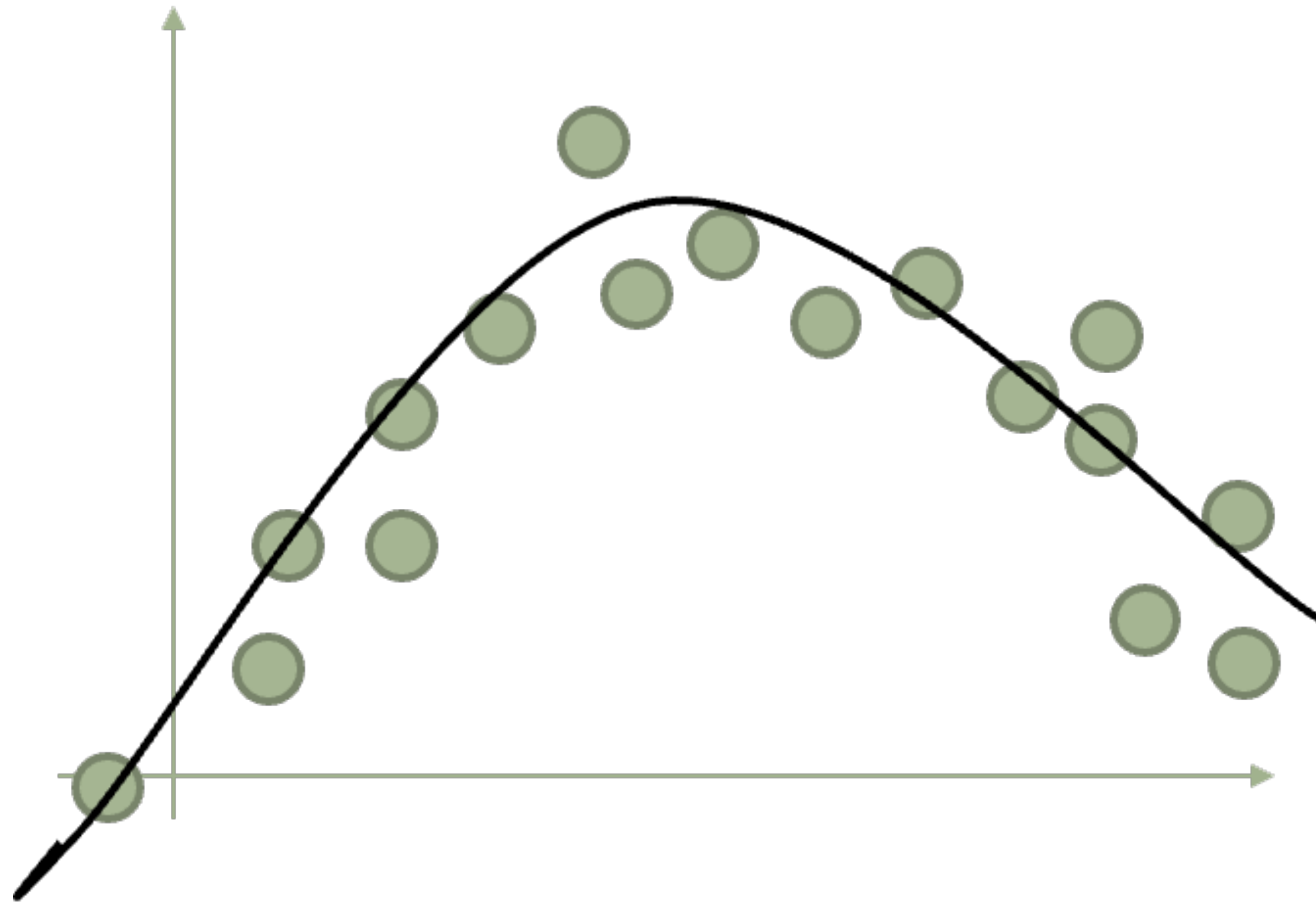
\$12 million or ~300 V100 GPUs for 3 months



Microsoft Research Blog. Oct 6, 2021.

Overfitting

Say we have the following graph

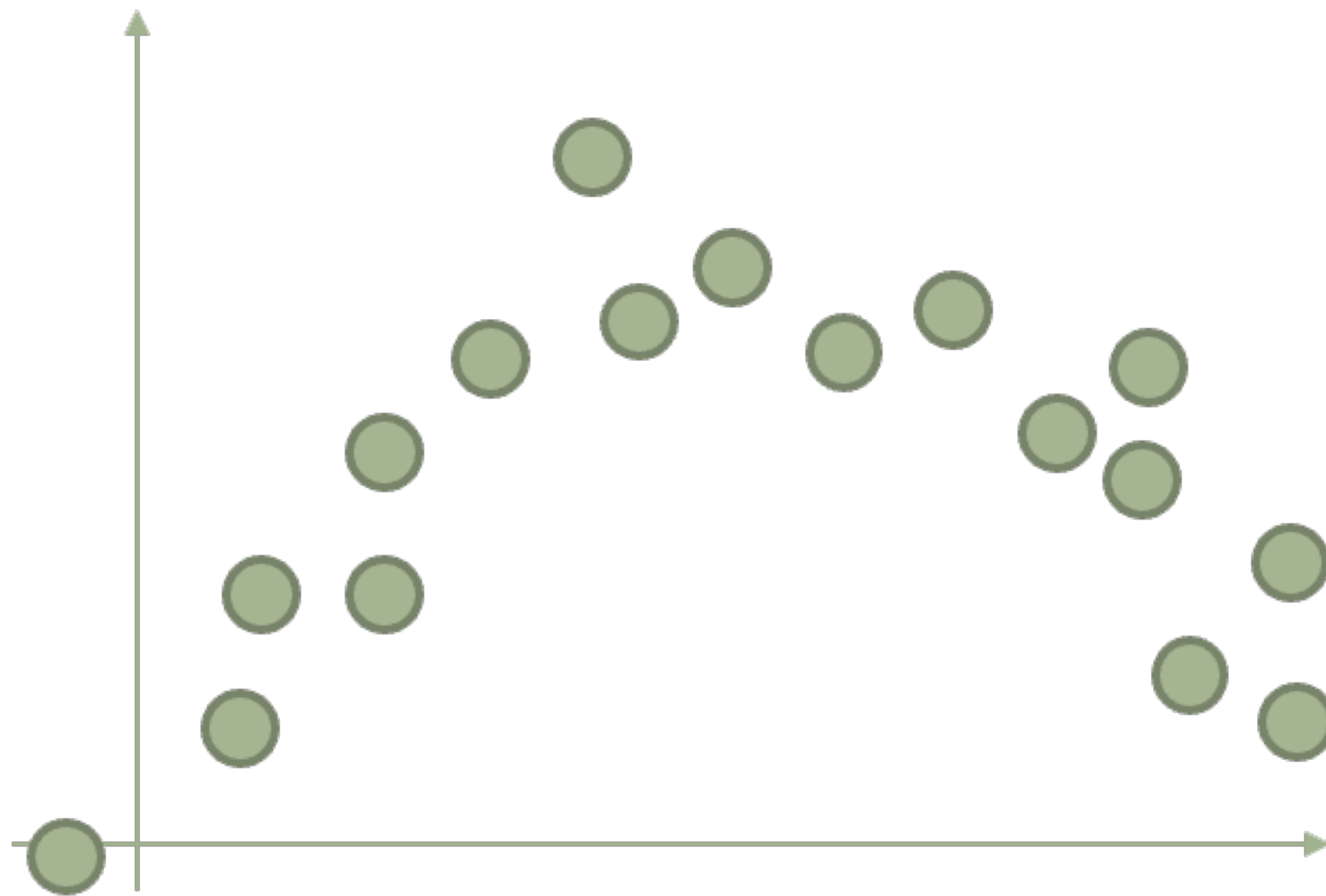


We know it should be a quadratic graph

Let's assume the best fit graph would be $y = -3x^2 + x + 0.5$

Overfitting

Say we have the following graph



We know it should be a quadratic graph

Let's assume the best fit graph would be $y = -3x^2 + x + 0.5$

If we fit with a higher order arbitrary polynomial function, $y = 0.4x^8 + 1.9x^7 - 1.4x^6 + \dots + 1.9$ -> overfit

Regularization

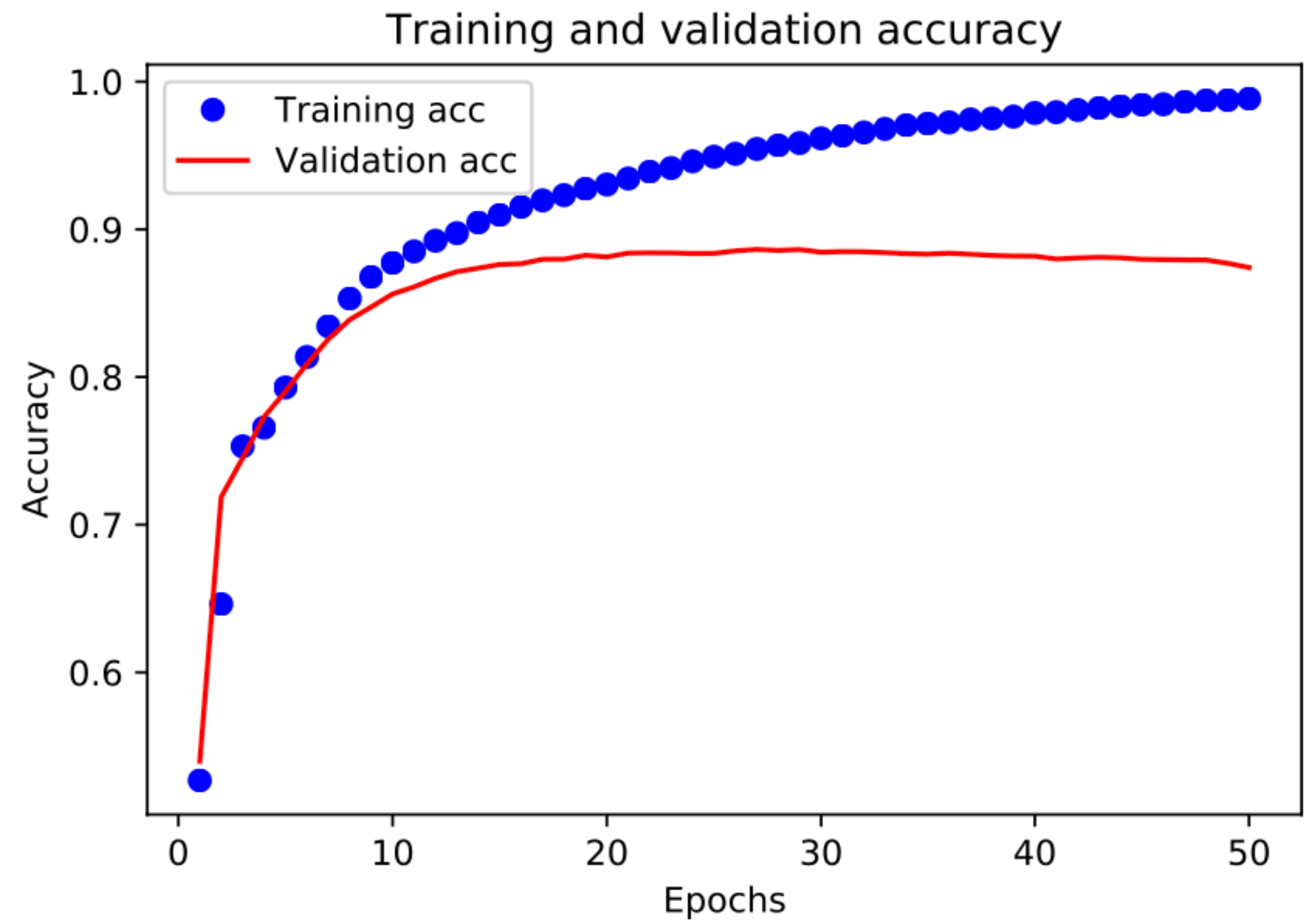
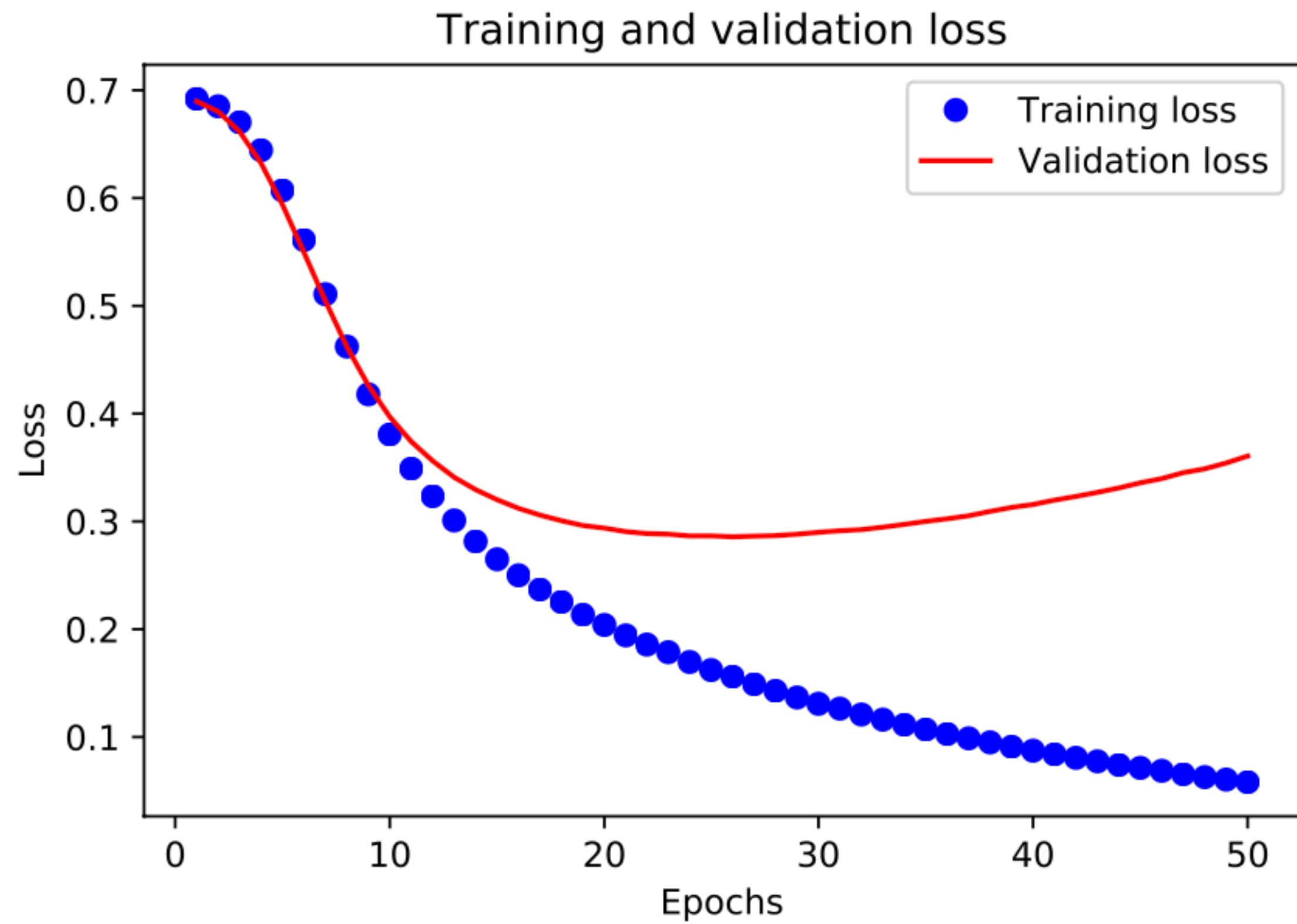
Weight Decay

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \underbrace{\frac{1}{N} \sum_{n=1}^N L(\mathbf{w}; x_n, y_n)}_{\text{empirical loss}} + \underbrace{\Omega(\mathbf{w})}_{\text{regularizer}}$$

$$\|\mathbf{w}\|_2 = \sqrt{\sum_{d=1}^D |w_d|^2} = \sqrt{\mathbf{w}^T \mathbf{w}}$$

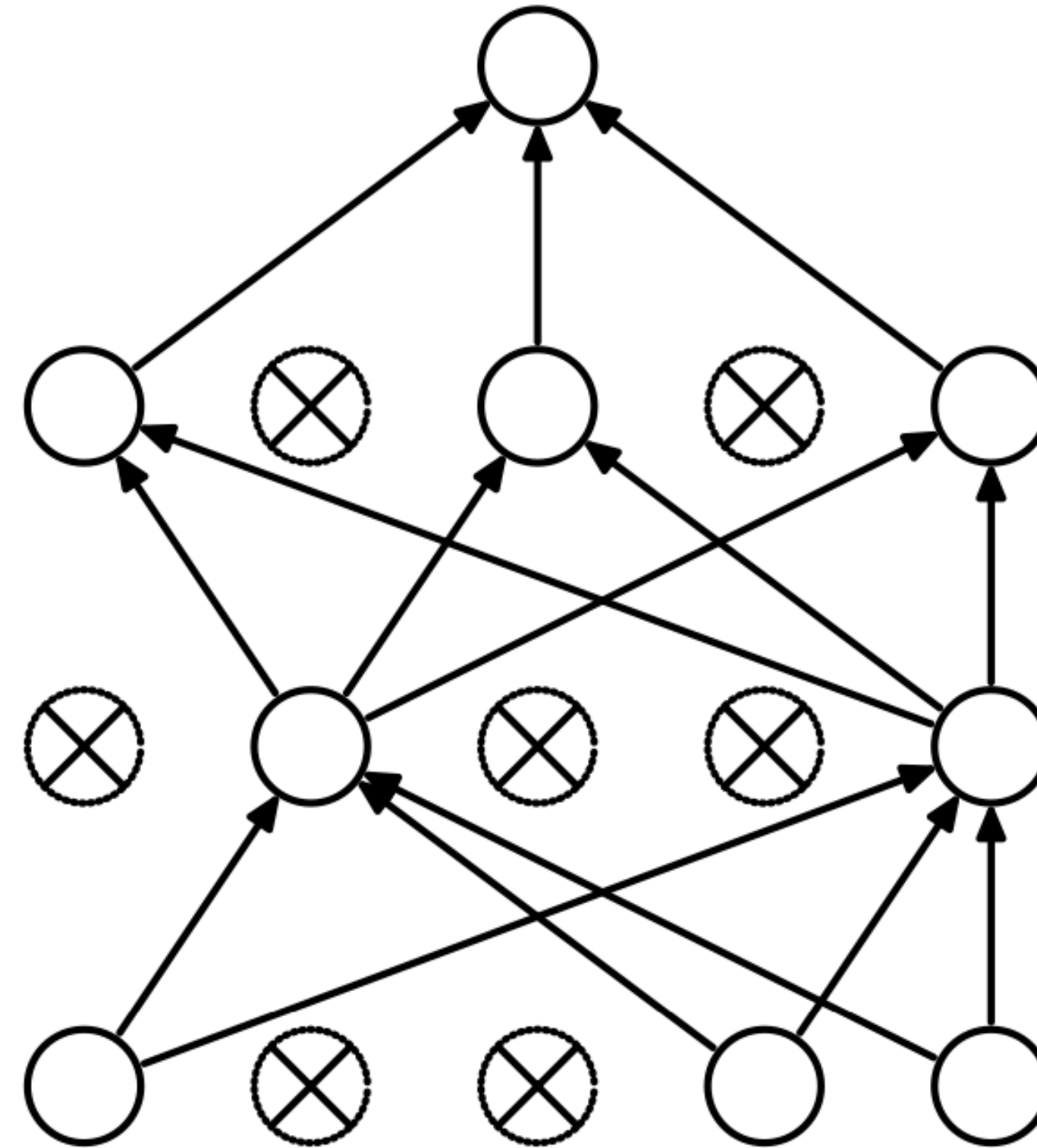
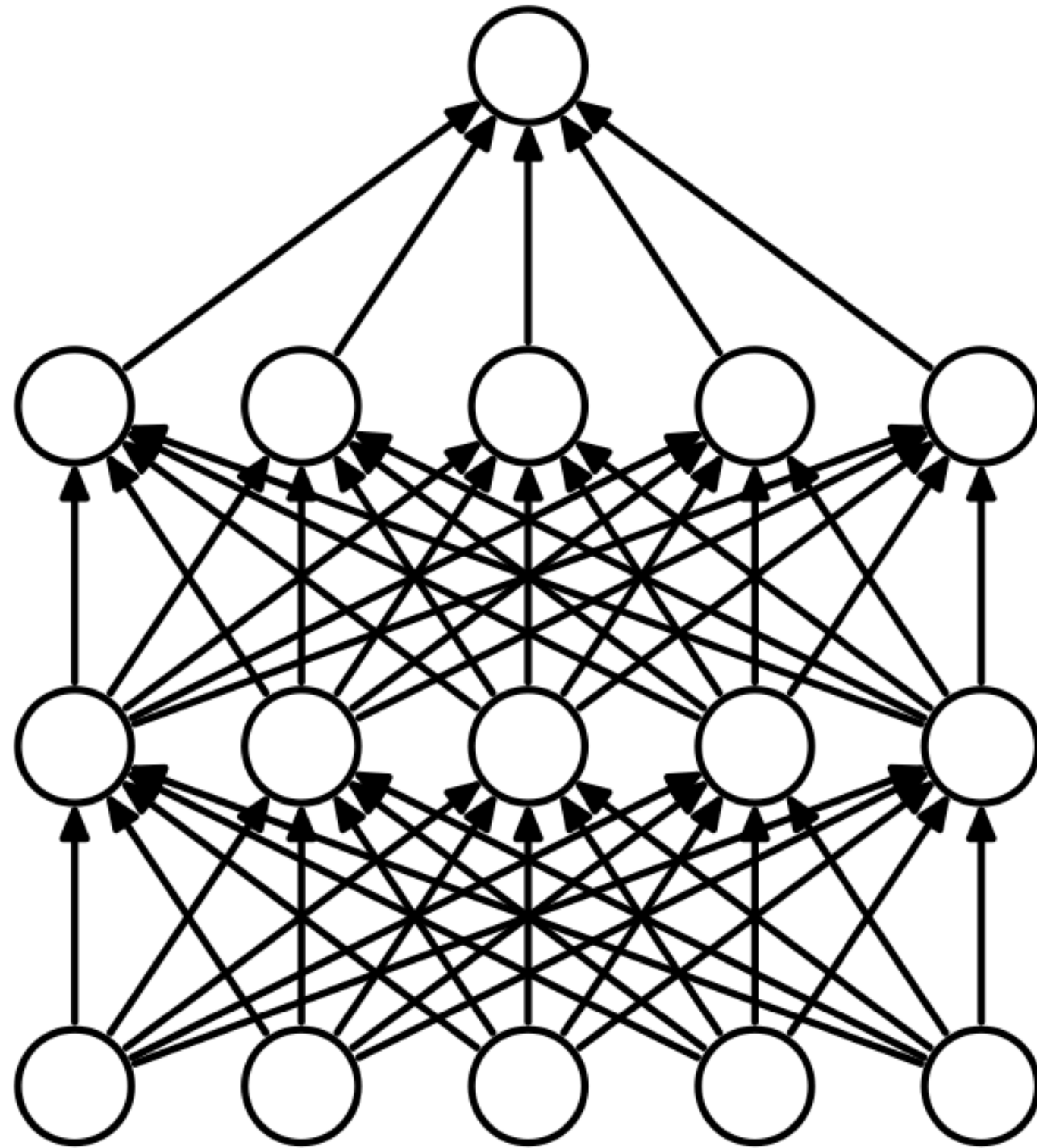
Regularization

Early stopping



Regularization

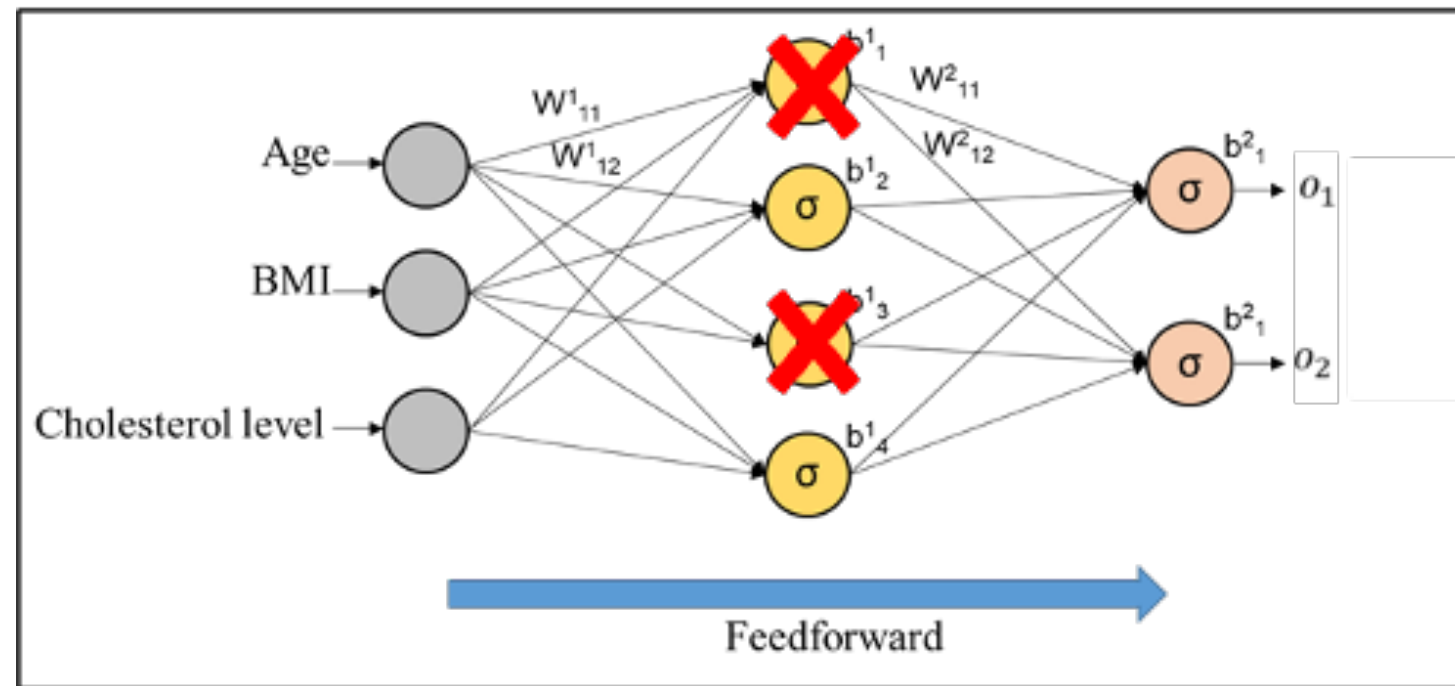
Dropout



Regularization

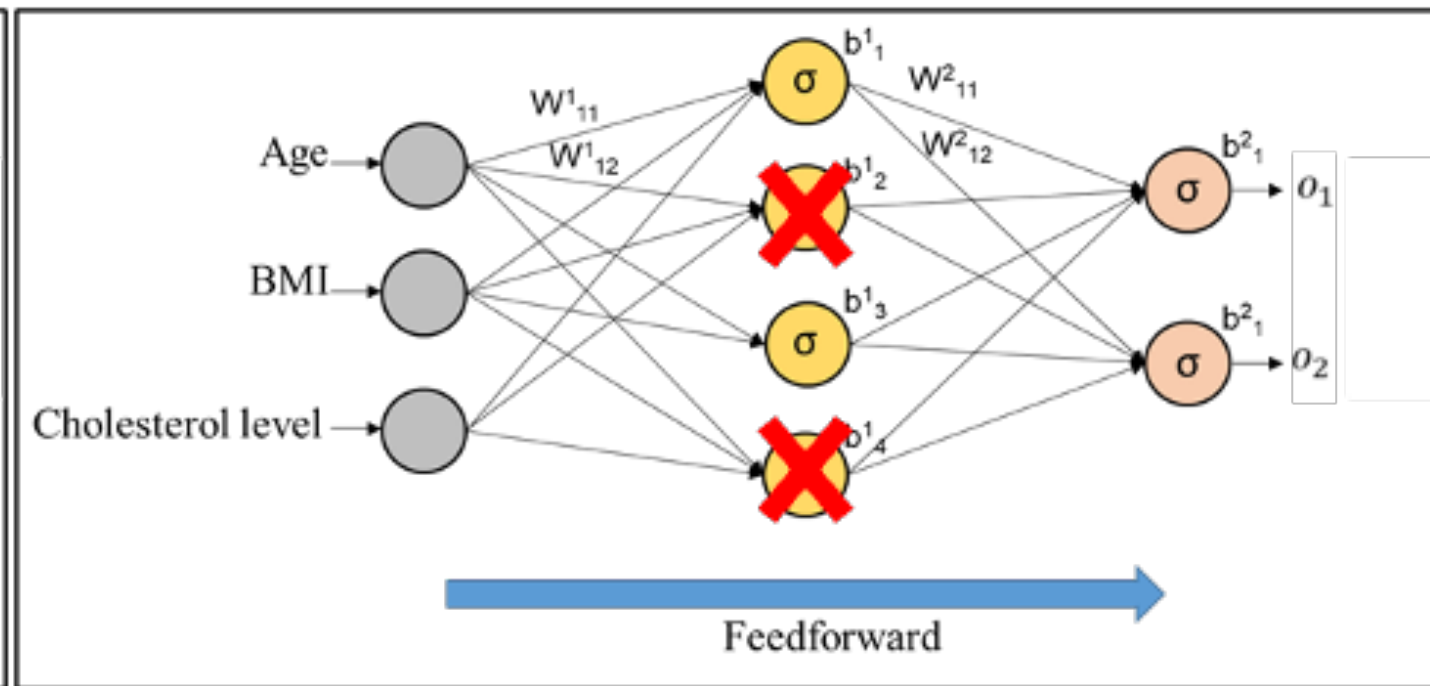
Dropout

1st iteration (Model1)



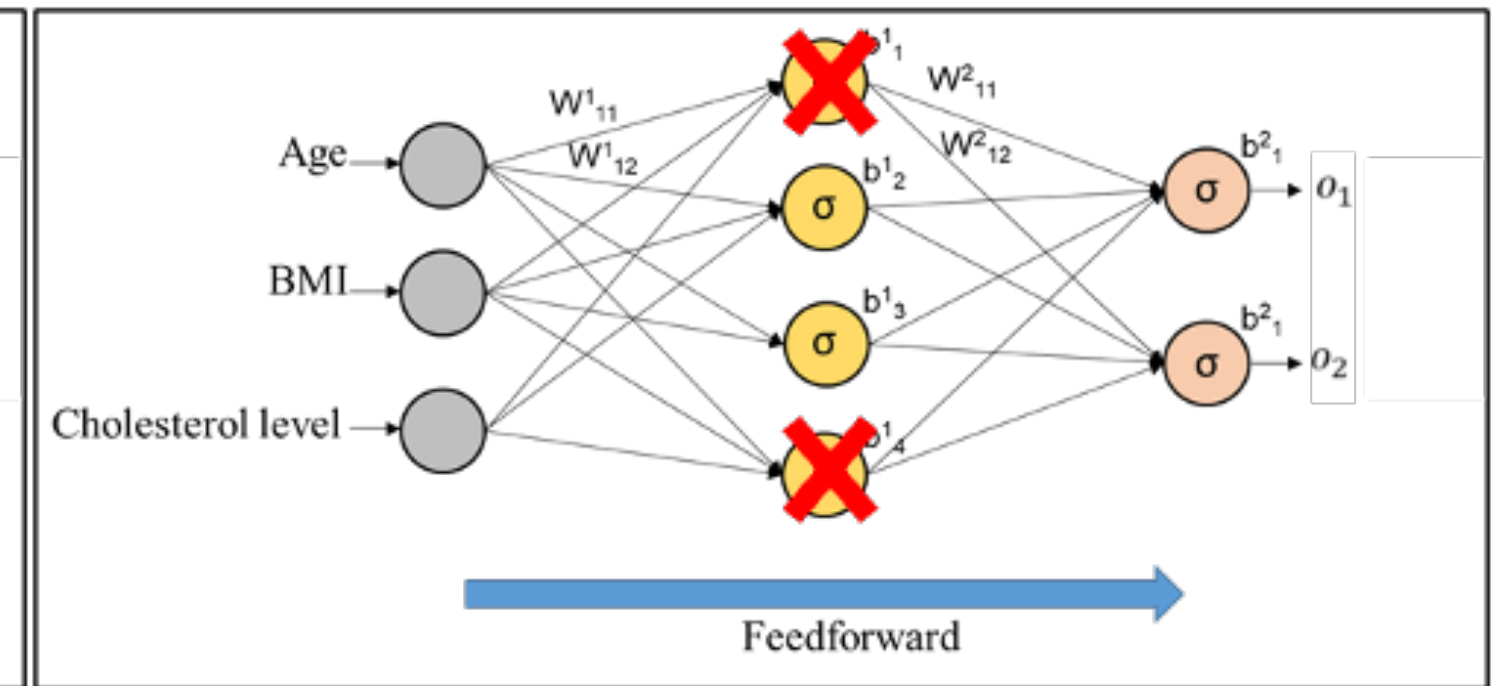
1. Forward
2. Compute gradient
3. Update weight

2nd iteration (Model2)



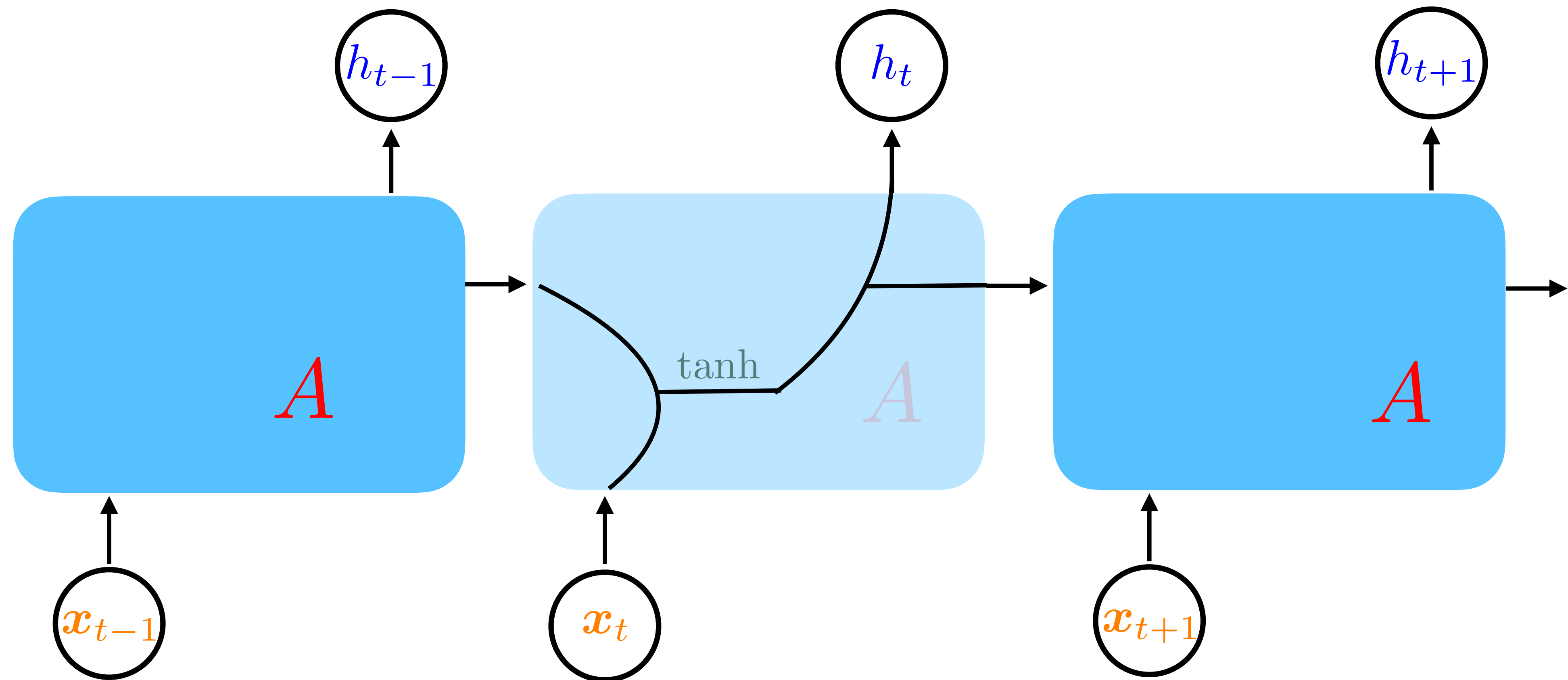
1. Forward
2. Compute gradient
3. Update weight

3rd iteration (Model3)



1. Forward
2. Compute gradient
3. Update weight

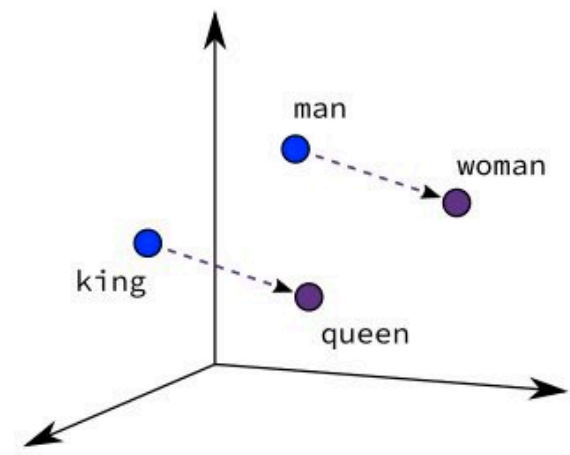
Recurrent Neural Networks (RNNs)



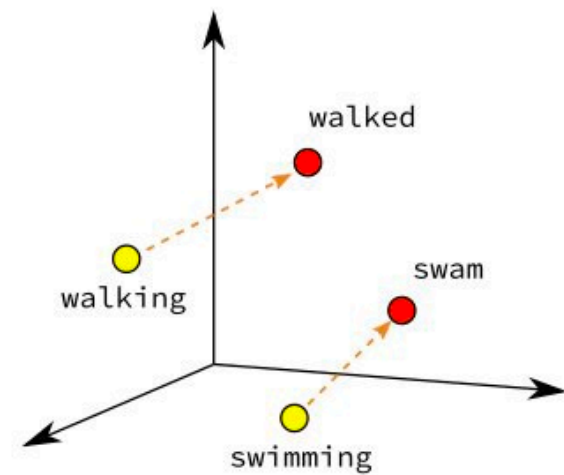
Word Embeddings



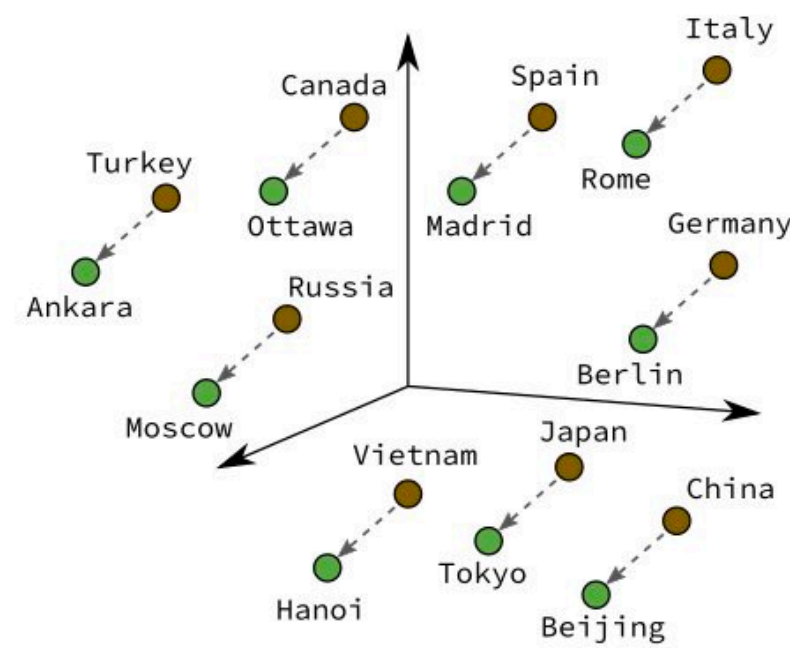
cat



Male-Female



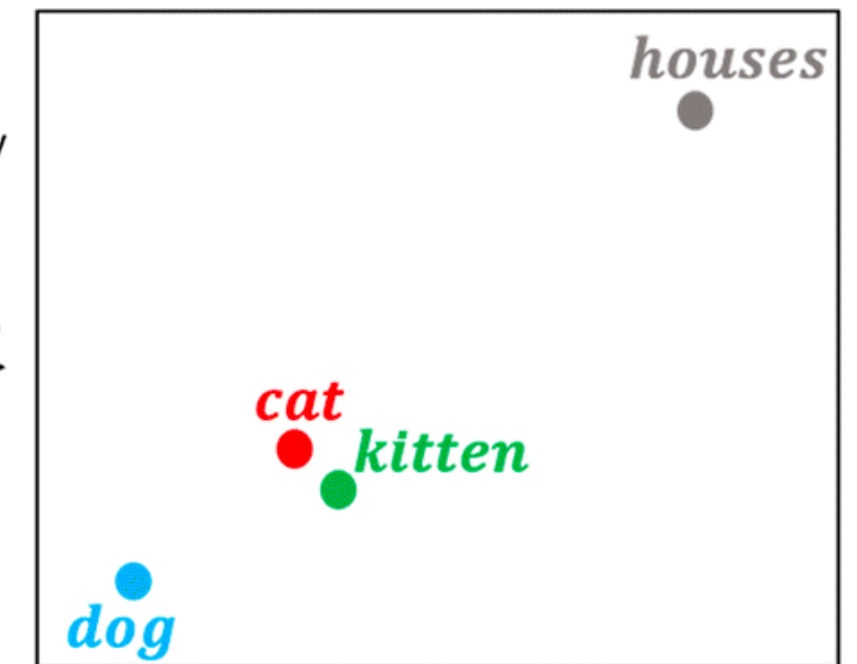
Verb Tense



Country-Capital

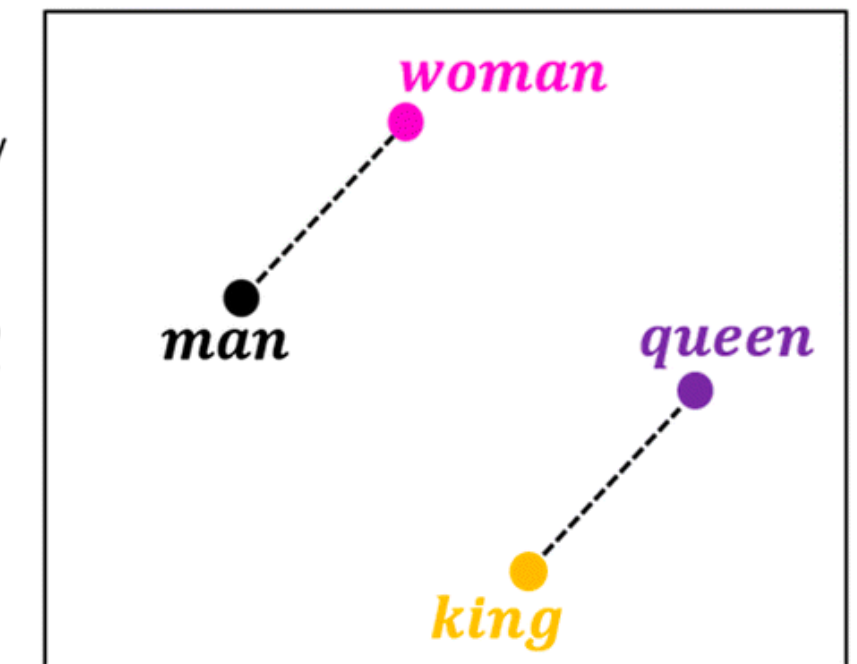
	living being	feline	human	gender	royalty	verb	plural
<i>cat</i>	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i>	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i>	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i>	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8

Dimensionality reduction of word embeddings from 7D to 2D



<i>man</i>	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i>	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i>	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i>	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

Dimensionality reduction of word embeddings from 7D to 2D



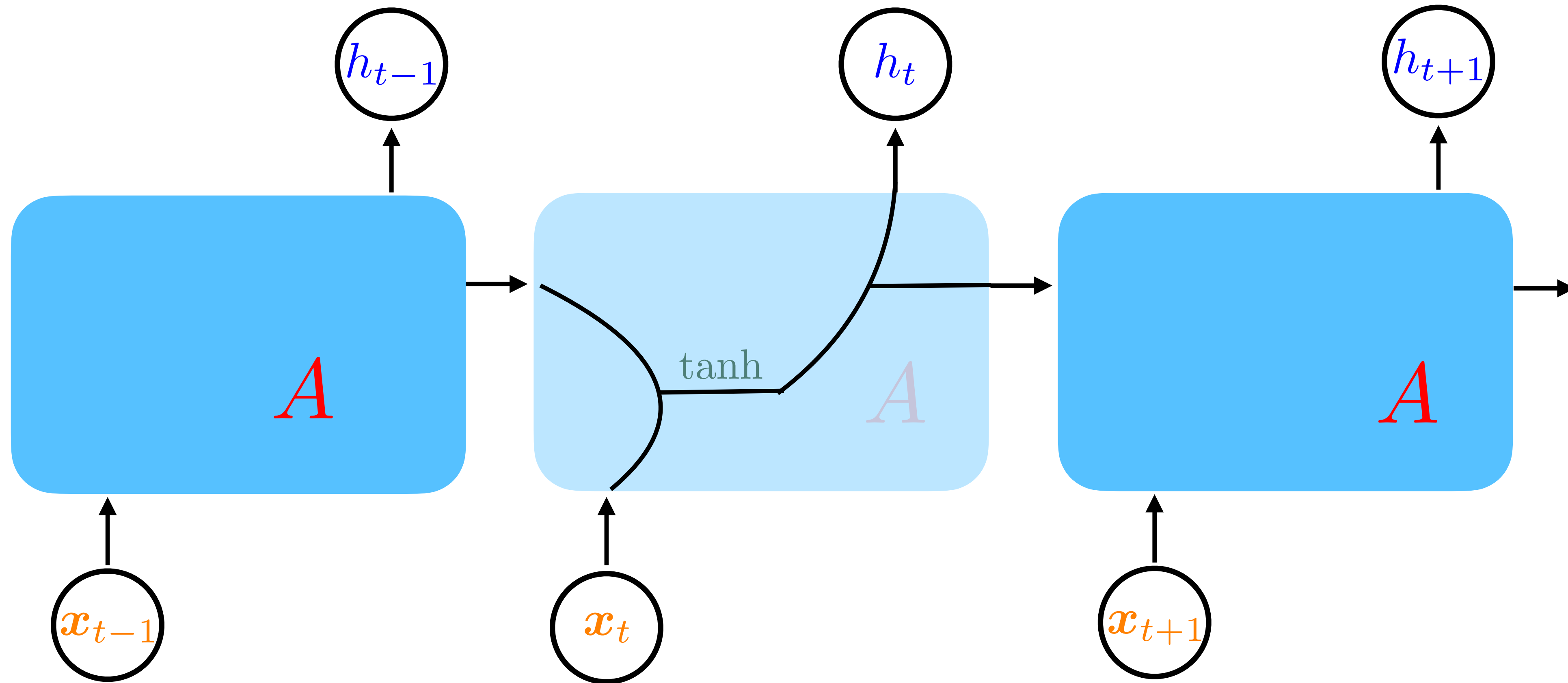
Word

Word embedding

Dimensionality reduction

Visualization of word embeddings in 2D

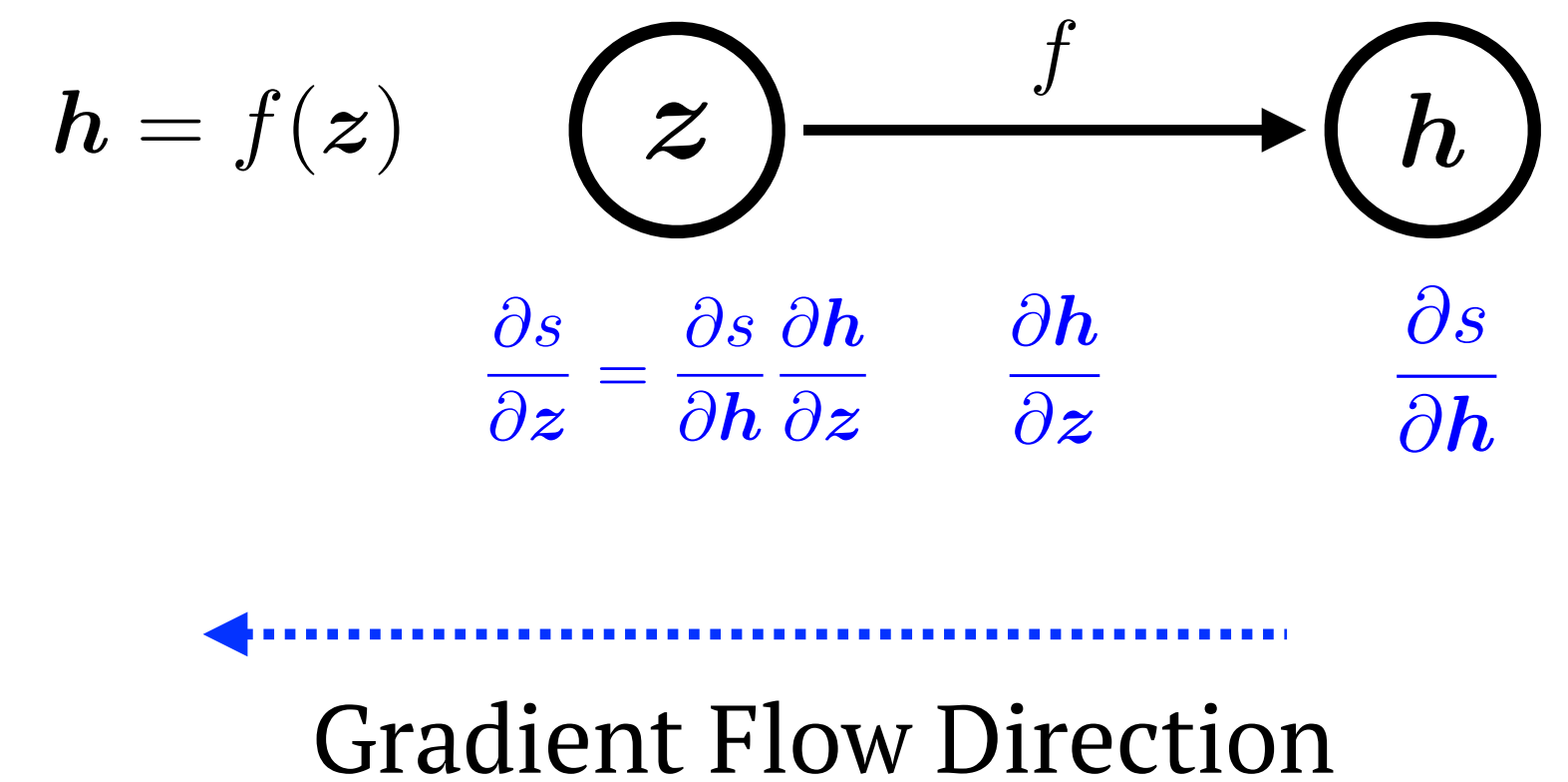
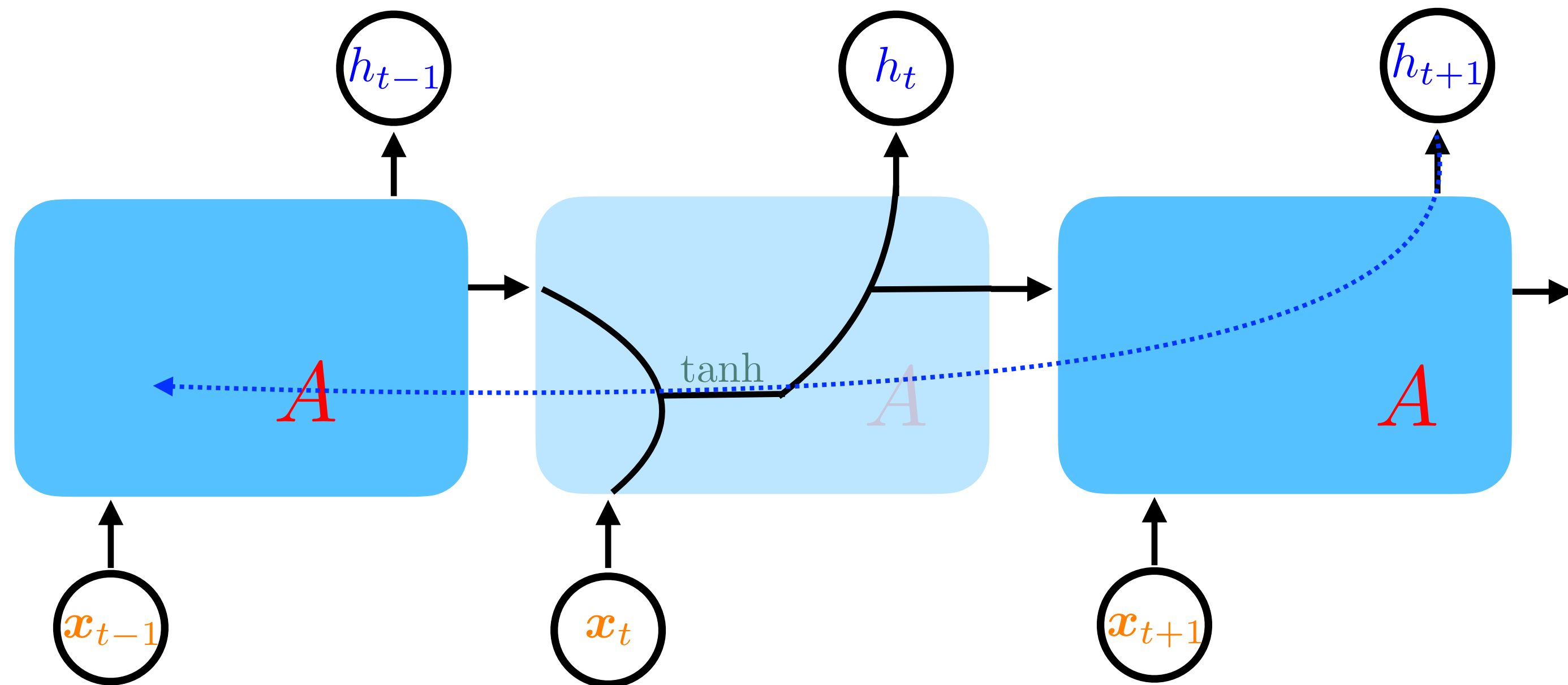
Recurrent Neural Networks (RNNs)



$$\mathbf{h}_t = \sigma(\mathbf{C}_x \mathbf{x}_t + \mathbf{C}_h \mathbf{h}_{t-1})$$

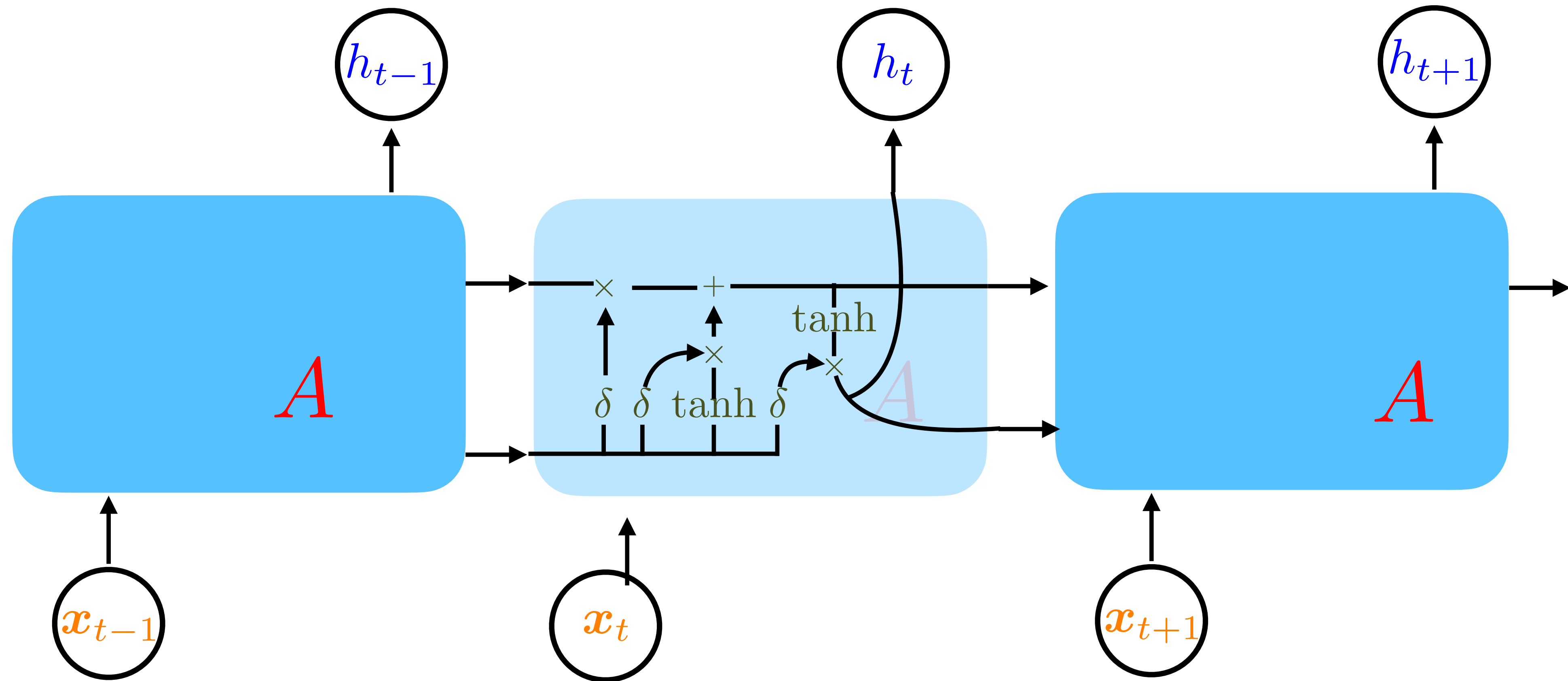
$$\mathbf{y}_t = f(\mathbf{W} \mathbf{h}_t + \mathbf{b})$$

Vanishing Gradient in RNNs



In general, the longer the path, the smaller the gradient signal.

Long Short-term Memory (LSTM)



Long Short-term Memory (LSTM)

$$\mathbf{i}_t = \sigma(\mathbf{I}_x \mathbf{x}_t + \mathbf{I}_h \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{f}_t = \sigma(\mathbf{F}_x \mathbf{x}_t + \mathbf{F}_h \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{o}_t = \sigma(\mathbf{O}_x \mathbf{x}_t + \mathbf{O}_h \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot g(\mathbf{C}_x \mathbf{x}_t + \mathbf{C}_h \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot g(\mathbf{c}_t)$$

$$\mathbf{y}_t = f(\mathbf{W} \mathbf{h}_t + \mathbf{b})$$

Part-of-Speech Tagging

INPUT:

Profits soared at Boeing Co., easily topping forecasts on Wall Street, as their CEO Alan Mulally announced first quarter results.

OUTPUT:

Profits/**N** soared/**V** at/**P** Boeing/**N** Co./**N** ,/, easily/**ADV** topping/**V** forecasts/**N** on/**P** Wall/**N** Street/**N** ,/, as/**P** their/**POSS** CEO/**N** Alan/**N** Mulally/**N** announced/**V** first/**ADJ** quarter/**N** results/**N** ./.

N = Noun

V = Verb

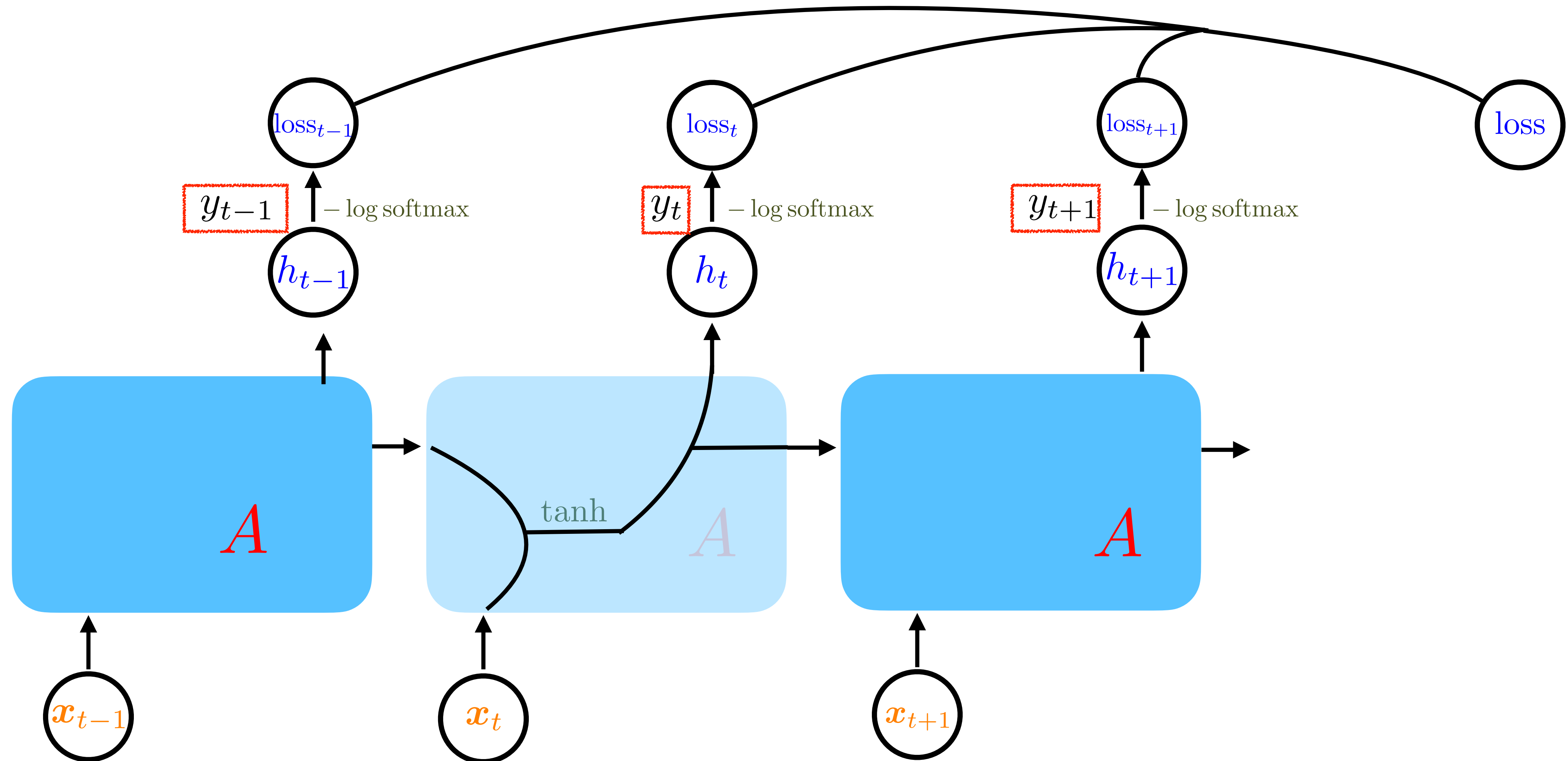
P = Preposition

Adv = Adverb

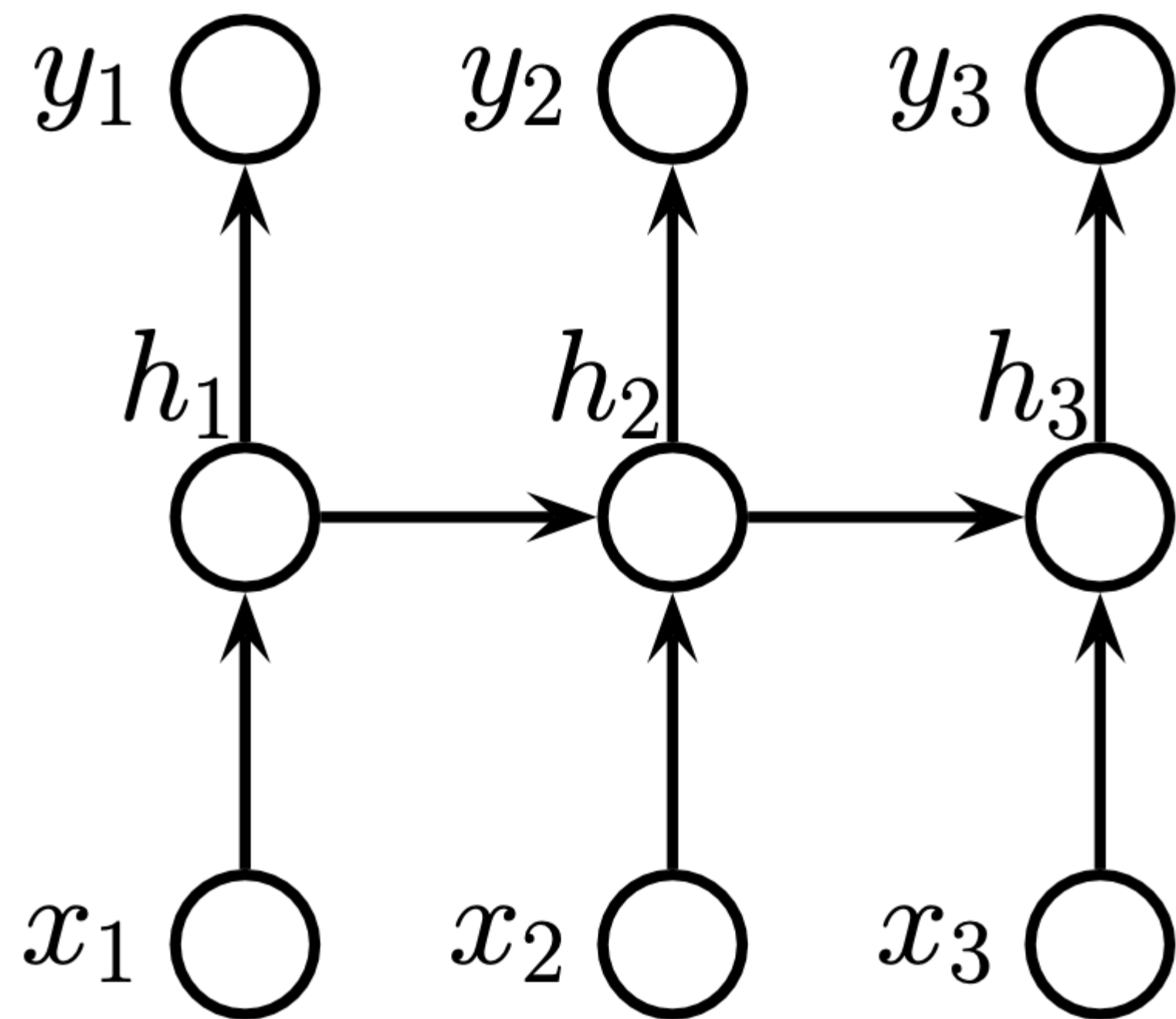
Adj = Adjective

...

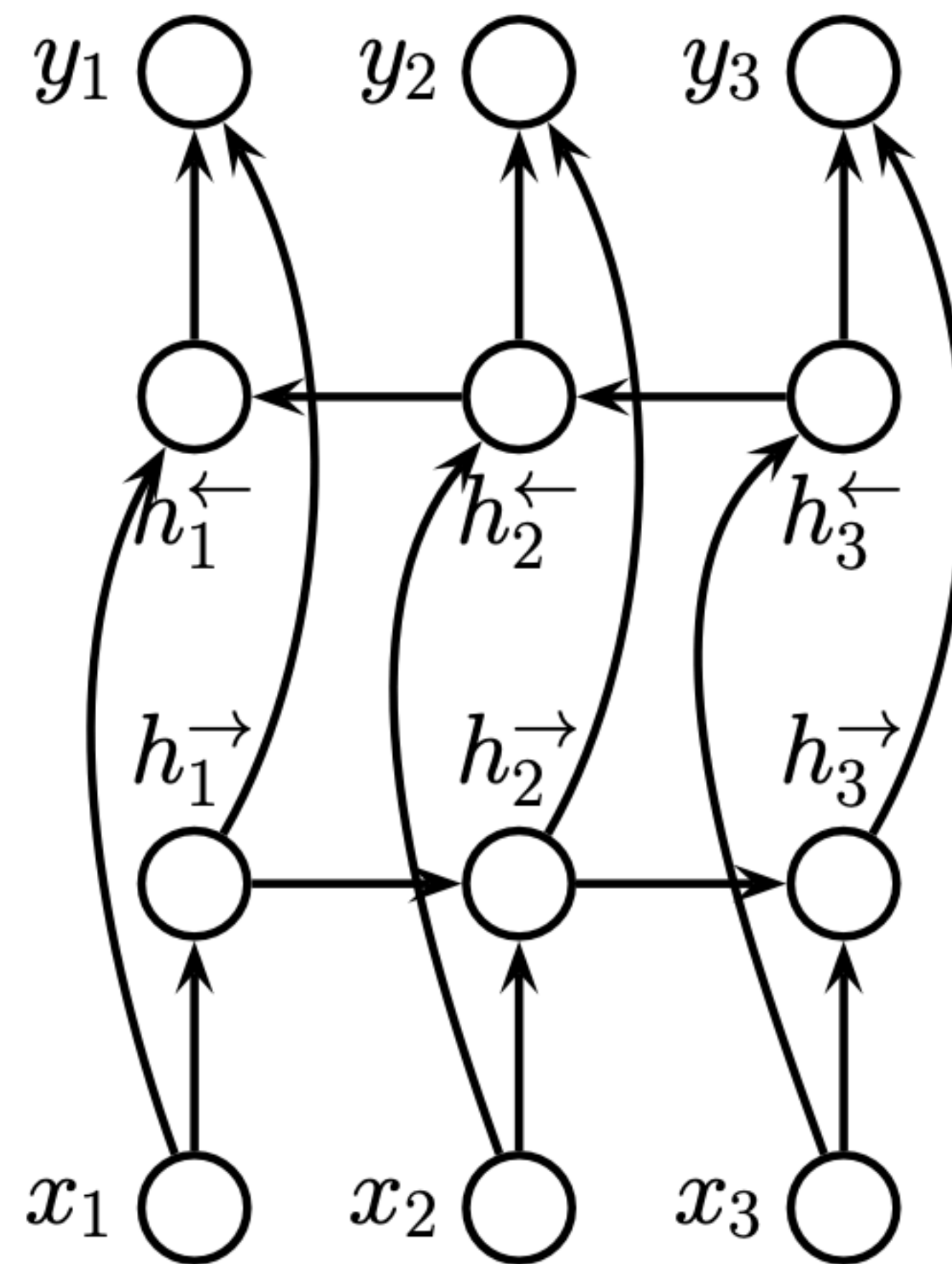
RNNs for Tagging



RNNs for Tagging



Left-to-right (Uni-directional)



Bidirectional