# Artificial Intelligence
# Techniques *in*
# Manufacturing
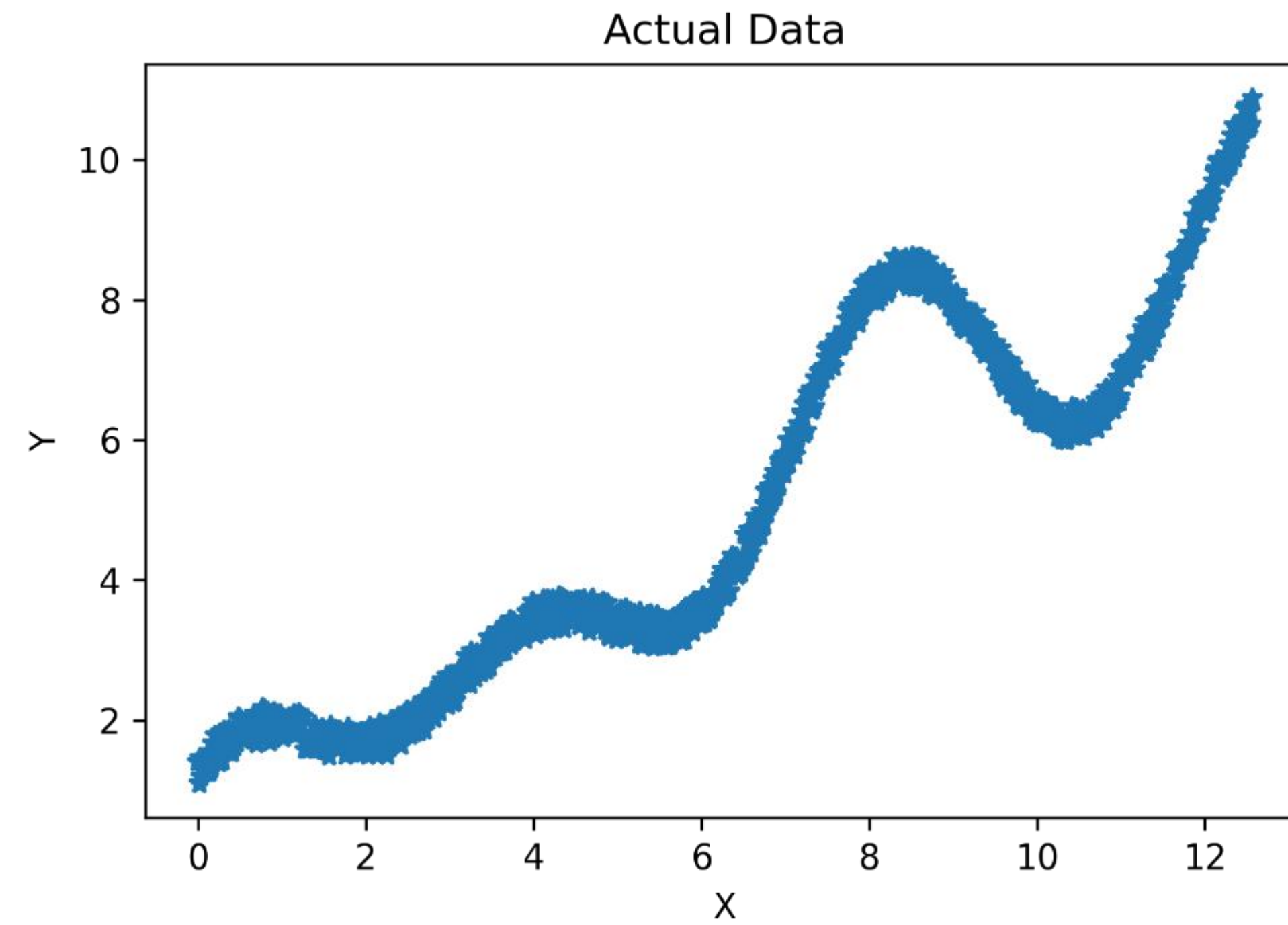
# Machine Learning: State of the Art

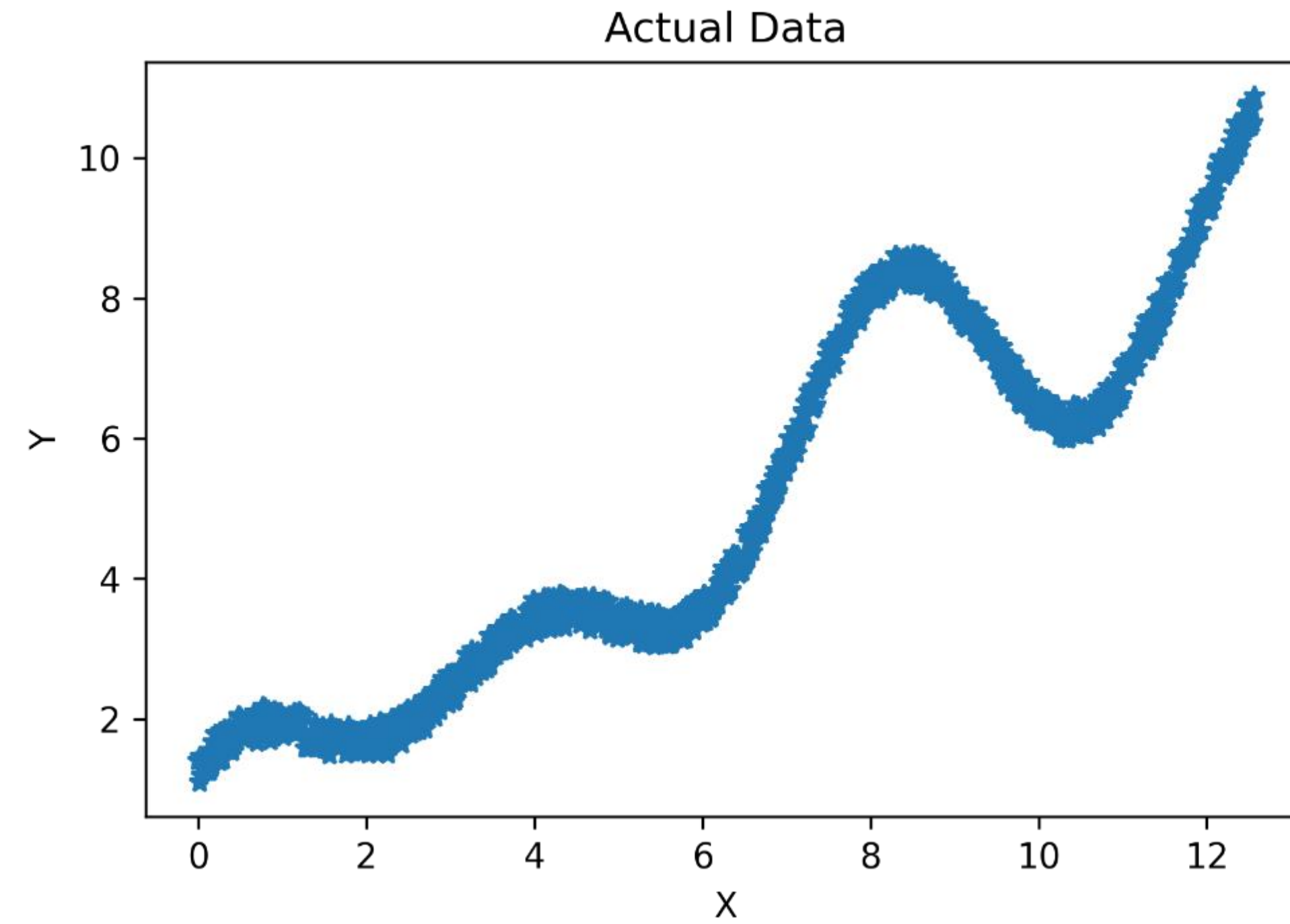| Data | Type | Examples |
|------|------|----------|
| Inhomogeneous Data | Gradient Boosting Machine | • XGBoost<br>• CatBoost<br>• LightGBM |
| Homogeneous Data | Artificial Neural Network | • Deep Neural Network<br>• Convolutional Neural Network<br>• Long Short-Term Memory Neural Network<br>• Transformer Network |

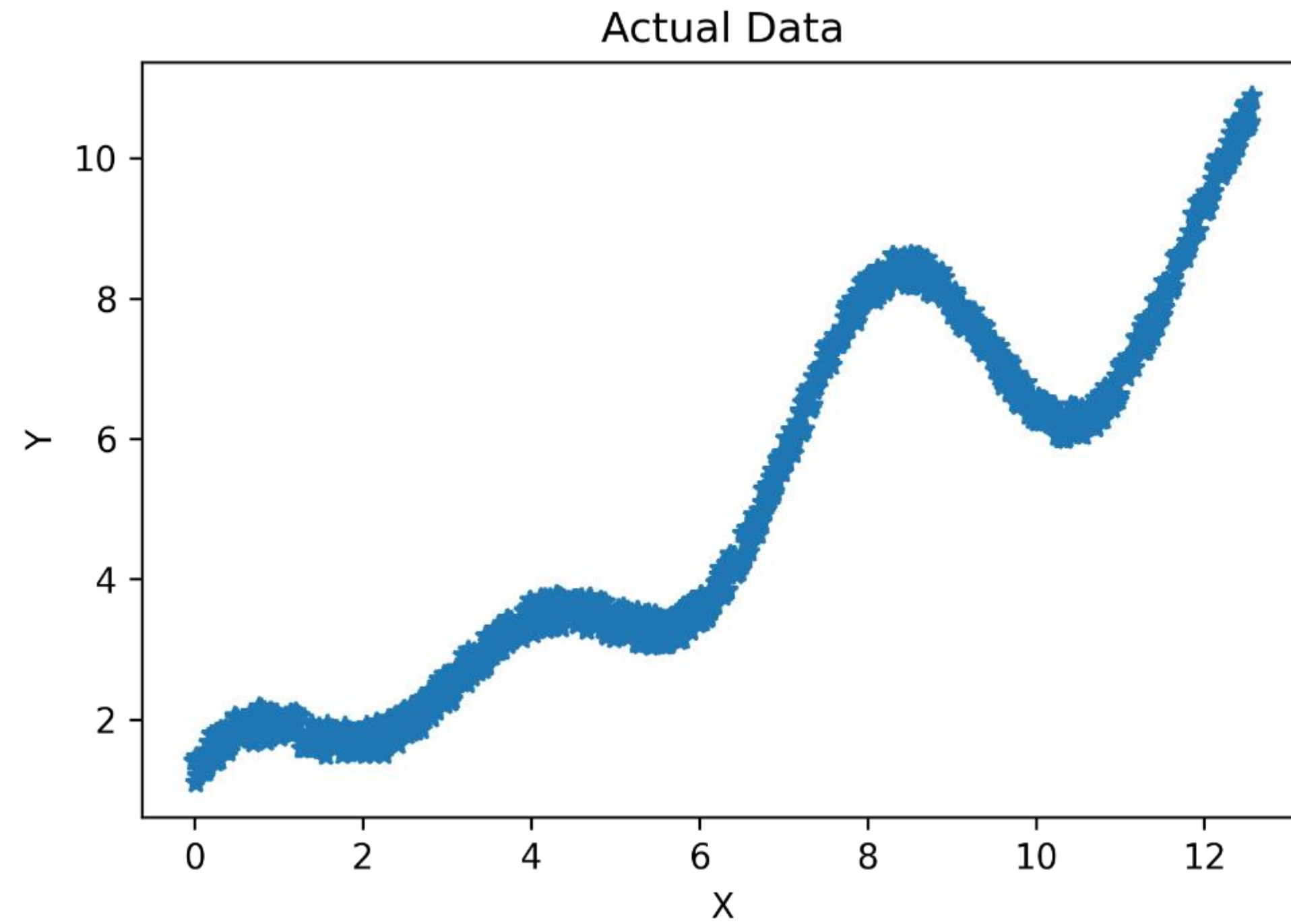# Artificial Neural Network

How do we fit a line to this data?

## Actual Data



$$y = f(x)$$

$$y = a_0 + a_1 x + a_2 x^2 + \cdots$$

Actual Data

$$y = f(x)$$

$$y = a_0 + a_1 x + a_2 x^2 + \cdots$$

$$y = a_0 + a_1 \sin(x) + a_2 \sin(2x) + \cdots$$

Actual Data

$$y = f(x)$$

$$y = a_0 + a_1 x + a_2 x^2 + \cdots$$

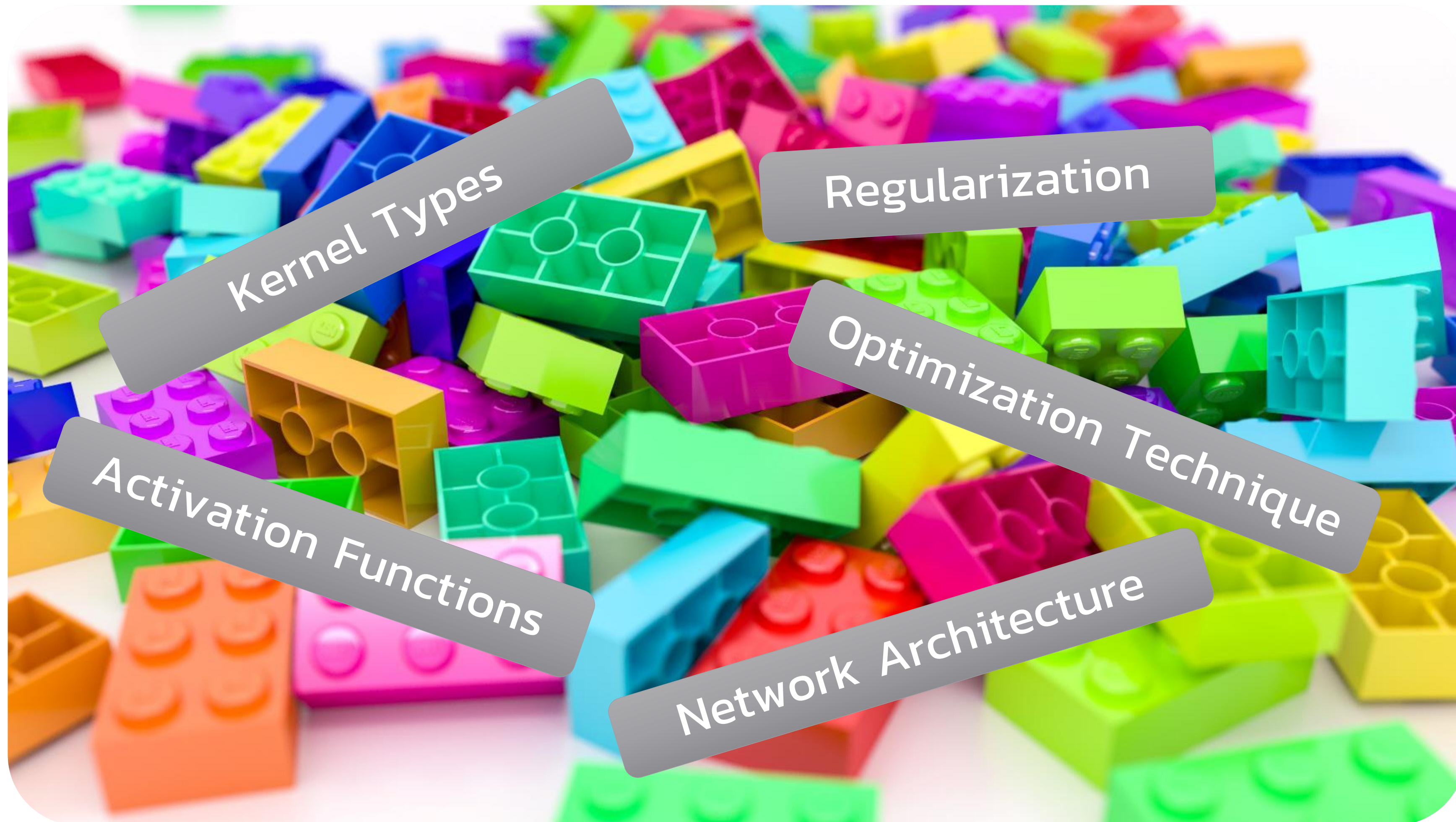$$y = a_0 + a_1 \sin(x) + a_2 \sin(2x) + \cdots$$

**Artificial Neural Network**

*Universal Approximator*

# Artificial Neural Network

# Artificial Neural Network

Kernel Types

Regularization

Optimization Technique

Activation Functions

Network Architecture
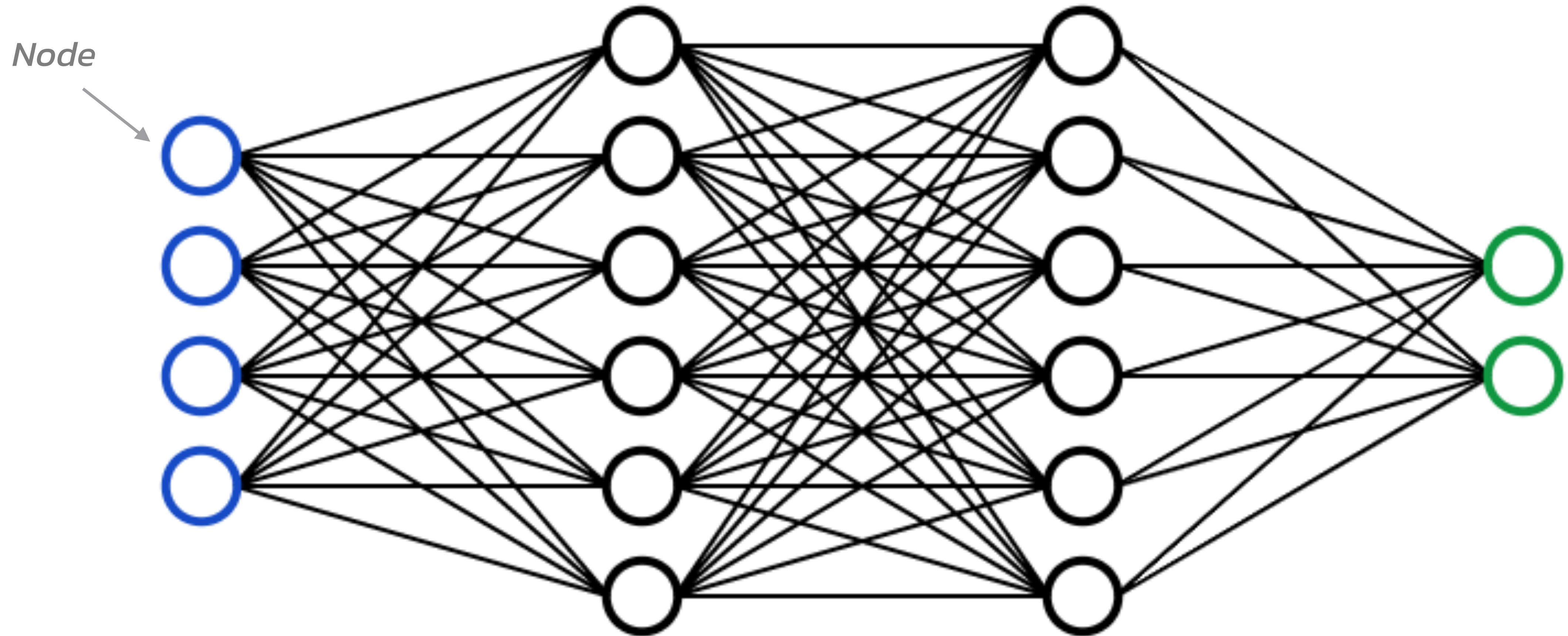
# Connection to Biological Neural Networks



This connection is not relevant nowadays.

# Architecture

Node

Input Layer

1ˢᵗ Hidden Layer

2ⁿᵈ Hidden Layer

Output Layer

# Hidden Node



$w_1, w_2, w_3$    $b$

Weights    Bias

$x_1$

$x_2$

$x_3$

$\Sigma$

$z$

$\phi$

$\phi$

$\phi$

$\phi$

Net Input Function

Activation Function

$z = X \cdot W + b$

$\phi = \phi(z)$

Input

Output

# Compared with Perceptron



$$\hat{y} = \begin{cases} 1, & if\ \phi(z) \geq 0 \\ -1, & Otherwise \end{cases}$$

Net input function

Activation function

Quantizer

Output

$$z = X \cdot W$$

$$\phi(z) = z$$

# Input Node

$x$

Single Input

$w = 1$
$b = 0$

$\Sigma$

Net Input Function

$z = x$

$z$

$\phi$

Activation Function

$\phi = z$

$\phi = x$

Single Output

# Output Node

Linear

#Parameters: 2

# Input Layer

# Hidden Layer

# Output Layer



$$z^{in} = x$$

$$\phi^{h1} = z^{in}$$

$$z^{h1} = w^{h1}\phi^{h1} + b^{h1}$$

$$\phi^{out} = \frac{1}{1 + e^{-z^{h1}}}$$

**Sigmoid Function**

$$z^{out} = w^{out}\phi^{out} + b^{out}$$

$$\hat{y} = z^{out}$$

**Input Layer** **Hidden Layer** **Output Layer**

$x \rightarrow \Sigma \xrightarrow{z^{in}} \phi \xrightarrow{\phi^{h1}} \Sigma \xrightarrow{z^{h1}} \phi \xrightarrow{\phi^{out}} \Sigma \xrightarrow{z^{out}} \phi \rightarrow \hat{y}$

$w^{h1}, b^{h1}$

$w^{out}, b^{out}$

$$\hat{y} = w^{out}\left[\frac{1}{1 + e^{-(w^{h1}x + b^{h1})}}\right] + b^{out}$$

S1

#Parameters: 4

$$\hat{y} = \Sigma_{i=1}^{10} \left[ w_i^{out} \frac{1}{1 + e^{-(w_i^{h1}x + b_i^{h1})}} \right] + b^{out}$$

S3

#Parameters: 4,353

# Deep Neural Network

$x$ — In — 256 — 256 — 128 — 64 — 32 — 16 — Out → $\hat{y}$

S6

#Parameters: 110,081

# Prediction



- Network
  - 2 hidden layers
  - Sigmoid activation
- Initialized weights and biases →
- Observation
  - $x = 1$
  - $y = 10$

| Variable | Init Val |
|----------|----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

| Variable | Init. Val |
|----------|-----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$x$

$1$

$z^{in}$

$1$

$\phi$

$\Sigma$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$\Sigma$

$z^{in}$

$\phi$

$w_1^{h2}, b_1^{h2}$

$\Sigma$

$z_1^{h2}$

$\phi$

$\phi^{h2}$

$w_1^{out}, w_2^{out}, b^{out}$

$\phi_1^{out}$

$\Sigma$

$z^{out}$

$\phi$

$\hat{y}$

$w_2^{h2}, b_2^{h2}$

$\phi^{h2}$

$\Sigma$

$z_2^{h2}$

$\phi$

$\phi_2^{out}$

| Variable | Init. Val |
|----------|-----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$z^{in} = x$$

| Variable | Init. Val |
|---|---|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$\phi^{h1} = z^{in}$$

$x$

$z^{in}$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$z^{in}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\phi^{h2}$

$w_1^{out}, w_2^{out}, b^{out}$

$z^{out}$

$\phi_1^{out}$

$w_2^{h2}, b_2^{h2}$

$z_2^{h2}$

$\phi^{h2}$

$\phi_2^{out}$

$\hat{y}$

| Variable | Init. Val |
|----------|-----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$z^{in} = w^{h1}\phi^{h1} + b^{h1}$$

$z^{in} = 1 \times 1 + 0 = 1$

$x$

$z^{in}$

$\phi$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$z^{in}$

$\phi$

$0.731$

$\phi^{h2}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\phi$

$0.731$

$\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$z_2^{h2}$

$\phi$

$w_1^{out}, w_2^{out}, b^{out}$

$\phi_1^{out}$

$\Sigma$

$z^{out}$

$\phi$

$\phi_2^{out}$

$\hat{y}$

| Variable | Init. Val |
|----------|-----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$\phi^{h2} = \frac{1}{1 + e^{-z^{in}}}$$

$$\phi^{h2} = \frac{1}{1 + e^{-1}} = 0.731$$

| Variable | Init. Val |
|---|---|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$x$

$z^{in}$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$z^{in}$

$\phi^{h2}$

0.731

$\phi^{h2}$

0.731

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

1.462

$w_2^{h2}, b_2^{h2}$

$z_2^{h2}$

2.193

$w_1^{out}, w_2^{out}, b^{out}$

$\phi_1^{out}$

$\phi_2^{out}$

$z^{out}$

$\hat{y}$

$$z_1^{h2} = w_1^{h2}\phi^{h2} + b_1^{h2}$$

$$z_2^{h2} = w_2^{h2}\phi^{h2} + b_2^{h2}$$

$z_1^{h2} = 2 \times 0.731 + 0 = 1.462$

$z_2^{h2} = 3 \times 0.731 + 0 = 2.193$

$x$

$z^{in}$

1

1

$\phi$

$\phi^{h1}$ 1

$w^{h1}, b^{h1}$

$z^{in}$ 1

$\phi$

0.731 $\phi^{h2}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\Sigma$ 1.462 $\phi$

0.811

$\phi_1^{out}$

$w_1^{out}, w_2^{out}, b^{out}$

$\Sigma$ $z^{out}$ $\phi$

$\hat{y}$

0.899

$\phi_2^{out}$

0.731 $\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$z_2^{h2}$

$\Sigma$ 2.193 $\phi$

| Variable | Init. Val |
|---|---|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$\phi_1^{out} = \frac{1}{1 + e^{-z_1^{h2}}}$$

$$\phi_2^{out} = \frac{1}{1 + e^{-z_2^{h2}}}$$

$x$

$z^{in}$

$\phi$

1

1

$\phi^{h1}$ 1

$w^{h1}, b^{h1}$

$z^{in}$

$\phi$

1

0.731

$\phi^{h2}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\Sigma$ $\phi$

1.462

0.811 $\phi_1^{out}$

$w_1^{out}, w_2^{out}, b^{out}$

$\Sigma$ $z^{out}$ $\phi$

7.745

0.731

$\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$z^{h2}$

$\Sigma$ $\phi$

2.193

0.899

$\phi_2^{out}$

$\hat{y}$

| Variable | Init. Val |
|----------|-----------|
| $w^{h1}$ | 1 |
| $w_1^{h2}$ | 2 |
| $w_2^{h2}$ | 3 |
| $w_1^{out}$ | 4 |
| $w_2^{out}$ | 5 |
| All biases | 0 |

$$z^{out} = w_1^{out}\phi_1^{out} + w_2^{out}\phi_2^{out} + b^{out}$$

# Loss

- Quantify how *"bad"* our model is.
- Mean Square Error (MSE)
  - $L = \frac{1}{N}\Sigma_{i=1}^{N}(y_i - \hat{y}_i)^2$
- In our example
  - $L = (10 - 7.745)^2 = 5.082$

# Parameter Update

- Backpropagation

  - Base on chain rule

- Computational graph

  - Compute input gradient signal at each node.

  - Send the gradient signal backward.

# Computational Graph



$$g^{out} = g^{in}\frac{\partial a^{out}}{\partial a^{in}}$$

# Computational Graph

$$g_1^{out} = g^{in} \frac{\partial a^{out}}{\partial a_1^{in}}$$

$$a_1^{in}$$

$$a^{out}$$

$$g_1^{out}$$
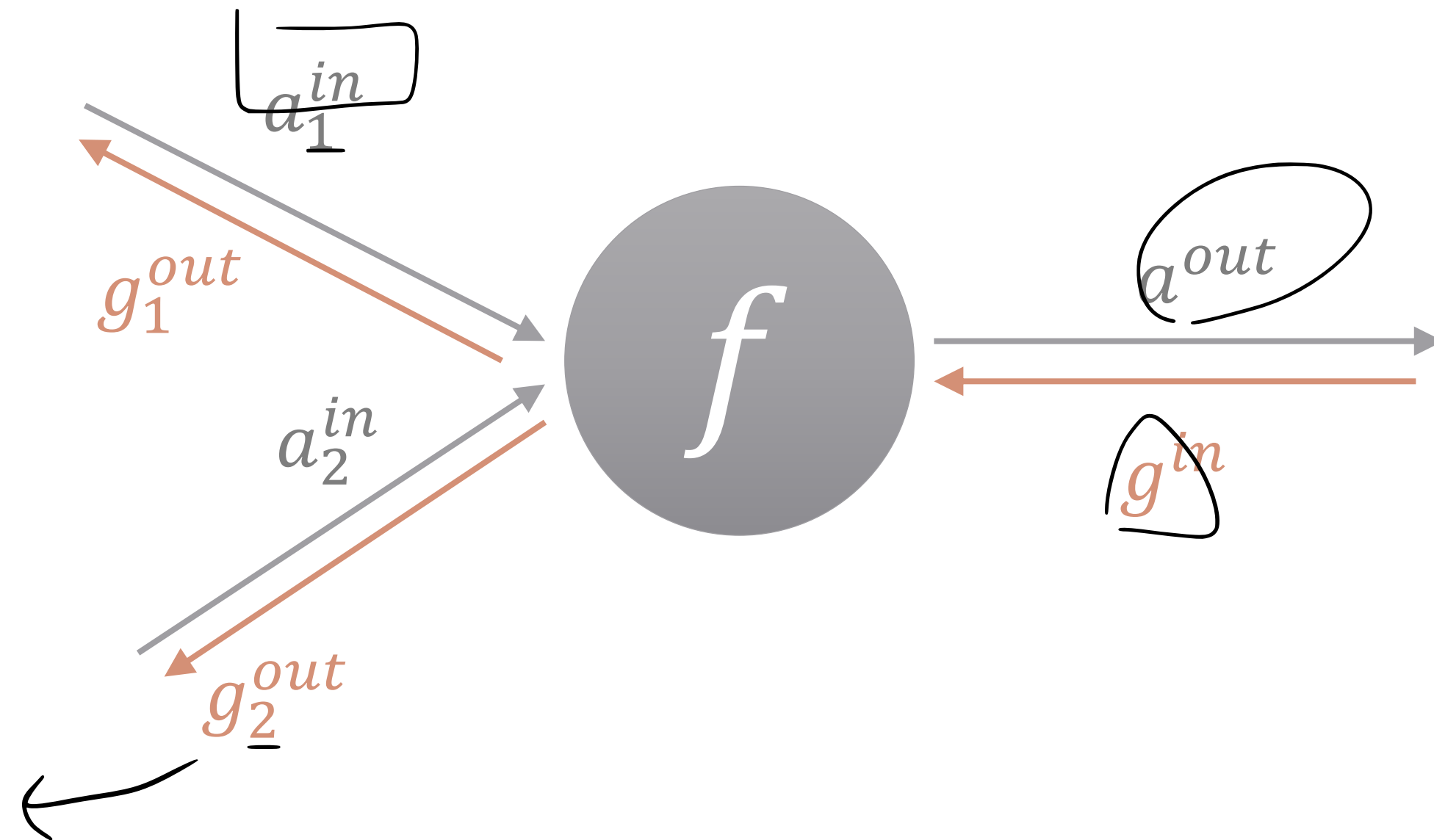
$$g^{in}$$

$$a_2^{in}$$

$$g_2^{out}$$

$f$
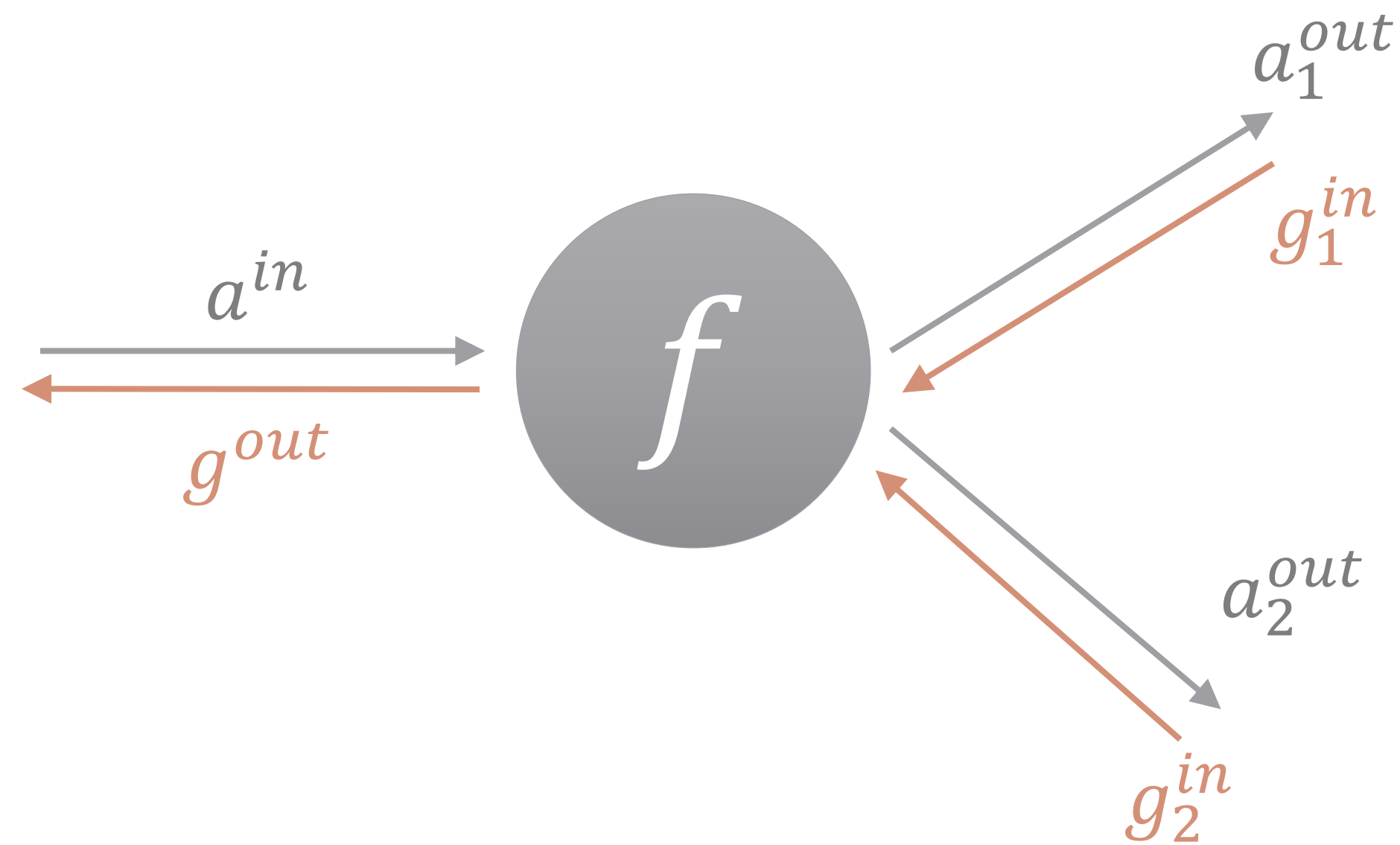
$$g_2^{out} = g^{in} \frac{\partial a^{out}}{\partial a_2^{in}}$$

# Computational Graph



$$g^{out} = g_1^{in} \frac{\partial a_1^{out}}{\partial a^{in}} + g_2^{in} \frac{\partial a_2^{out}}{\partial a^{in}}$$

$z^{in}$

$\Sigma$   $\phi$

$x$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$\Sigma$   $z^{in}$   $\phi$

$\phi^{h2}$

$w_1^{h2}, b_1^{h2}$

$\Sigma$   $z_1^{h2}$   $\phi$

$\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$\Sigma$   $z_2^{h1}$   $\phi$

$\phi_1^{out}$

$\phi_2^{out}$

$w_1^{out}, w_2^{out}, b^{out}$

$\Sigma$   $z^{out}$   $\phi$

$\hat{y}$

$L$

$1$

$L$

Add another $"L"$ node with gradient input of 1.

$z^{in}$

$\phi^{h1}$

$x$

$w^{h1}, b^{h1}$

$z^{in}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\phi^{h2}$

$w_1^{out}, w_2^{out}, b^{out}$

$z^{out}$

$\phi_1^{out}$

$\hat{y}$

$-4.508$

$L$

$w_2^{h2}, b_2^{h2}$

$z_2^{h1}$

$\phi^{h2}$

$\phi_2^{out}$

$1$

$L$

$$g^{in} = \frac{\partial L}{\partial L} = 1$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y) = -4.508$$

$$g^{out} = g^{in}\frac{\partial L}{\partial \hat{y}} = -4.508$$

$z^{in}$

$x$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$z^{in}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$z_2^{h1}$

$\phi^{h2}$

$\phi_1^{out}$

$\phi_2^{out}$

$w_1^{out}, w_2^{out}, b^{out}$

$z^{out}$

$-4.508$

$-4.508$

$\hat{y}$

$L$

1

$L$

$g^{in} = -4.508$

$\dfrac{\partial \hat{y}}{\partial z^{out}} = 1$

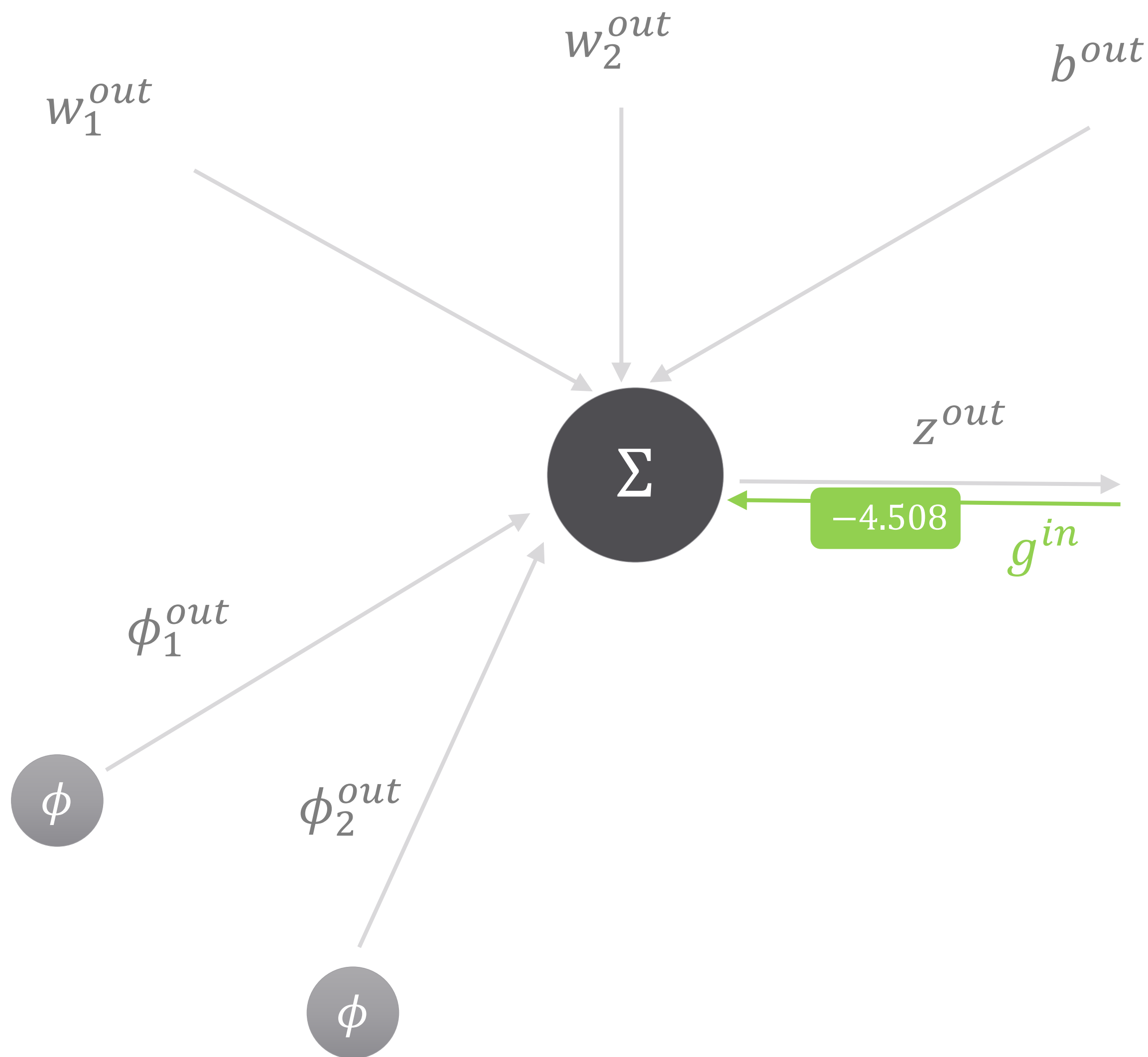$g^{out} = g^{in} \dfrac{\partial \hat{y}}{\partial z^{out}} = -4.508$

$w_2^{out}$

$w_1^{out}$

$b^{out}$

$z^{out} = w_1^{out}\phi_1^{out} + w_2^{out}\phi_2^{out} + b^{out}$

$g^{in} = -4.508$

$\Sigma$

$z^{out}$

$-4.508$

$g^{in}$

$\phi_1^{out}$

$\phi$

$\phi_2^{out}$

$\phi$

$$z^{out} = w_1^{out}\phi_1^{out} + w_2^{out}\phi_2^{out} + b^{out}$$

$$g^{in} = -4.508$$

$$\frac{\partial z^{out}}{\partial \phi_1^{out}} = w_1^{out} = 4$$

$$\frac{\partial z^{out}}{\partial w_1^{out}} = \phi_1^{out} = 0.811$$

$$\frac{\partial z^{out}}{\partial \phi_2^{out}} = w_2^{out} = 5$$

$$\frac{\partial z^{out}}{\partial w_2^{out}} = \phi_2^{out} = 0.899$$

$$\frac{\partial z^{out}}{\partial b^{out}} = 1$$

$$z^{out} = w_1^{out} \phi_1^{out} + w_2^{out} \phi_2^{out} + b^{out}$$

$$g^{in} = -4.508$$

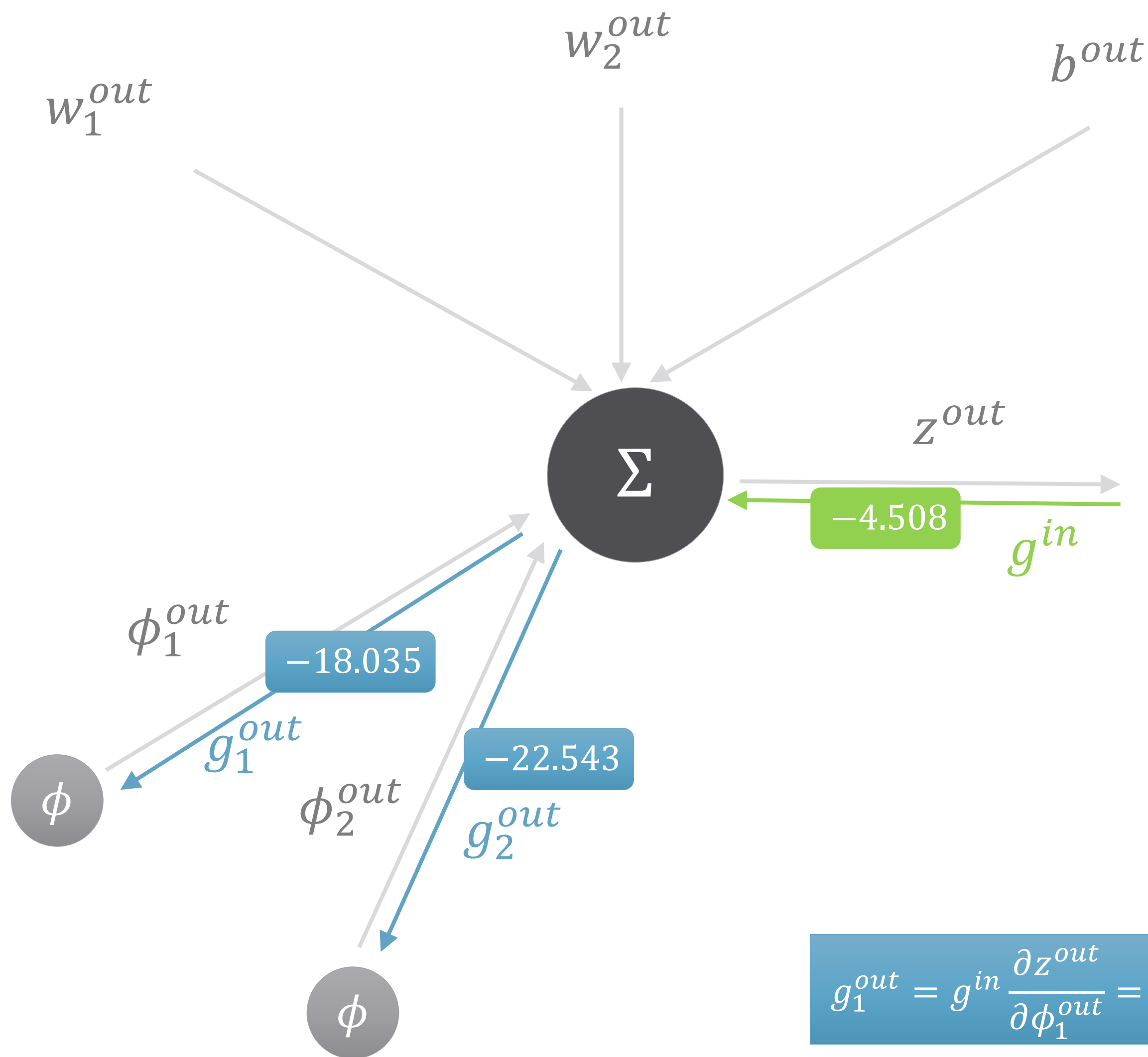$$\frac{\partial z^{out}}{\partial \phi_1^{out}} = w_1^{out} = 4$$

$$\frac{\partial z^{out}}{\partial w_1^{out}} = \phi_1^{out} = 0.811$$

$$\frac{\partial z^{out}}{\partial \phi_2^{out}} = w_2^{out} = 5$$

$$\frac{\partial z^{out}}{\partial w_2^{out}} = \phi_2^{out} = 0.899$$

$$\frac{\partial z^{out}}{\partial b^{out}} = 1$$

$$g_1^{out} = g^{in} \frac{\partial z^{out}}{\partial \phi_1^{out}} = -4.508 \times 4 = -18.035$$

$$g_2^{out} = g^{in} \frac{\partial z^{out}}{\partial \phi_2^{out}} = -4.508 \times 5 = -22.543$$

$$z^{out} = w_1^{out}\phi_1^{out} + w_2^{out}\phi_2^{out} + b^{out}$$

$$g^{in} = -4.508$$

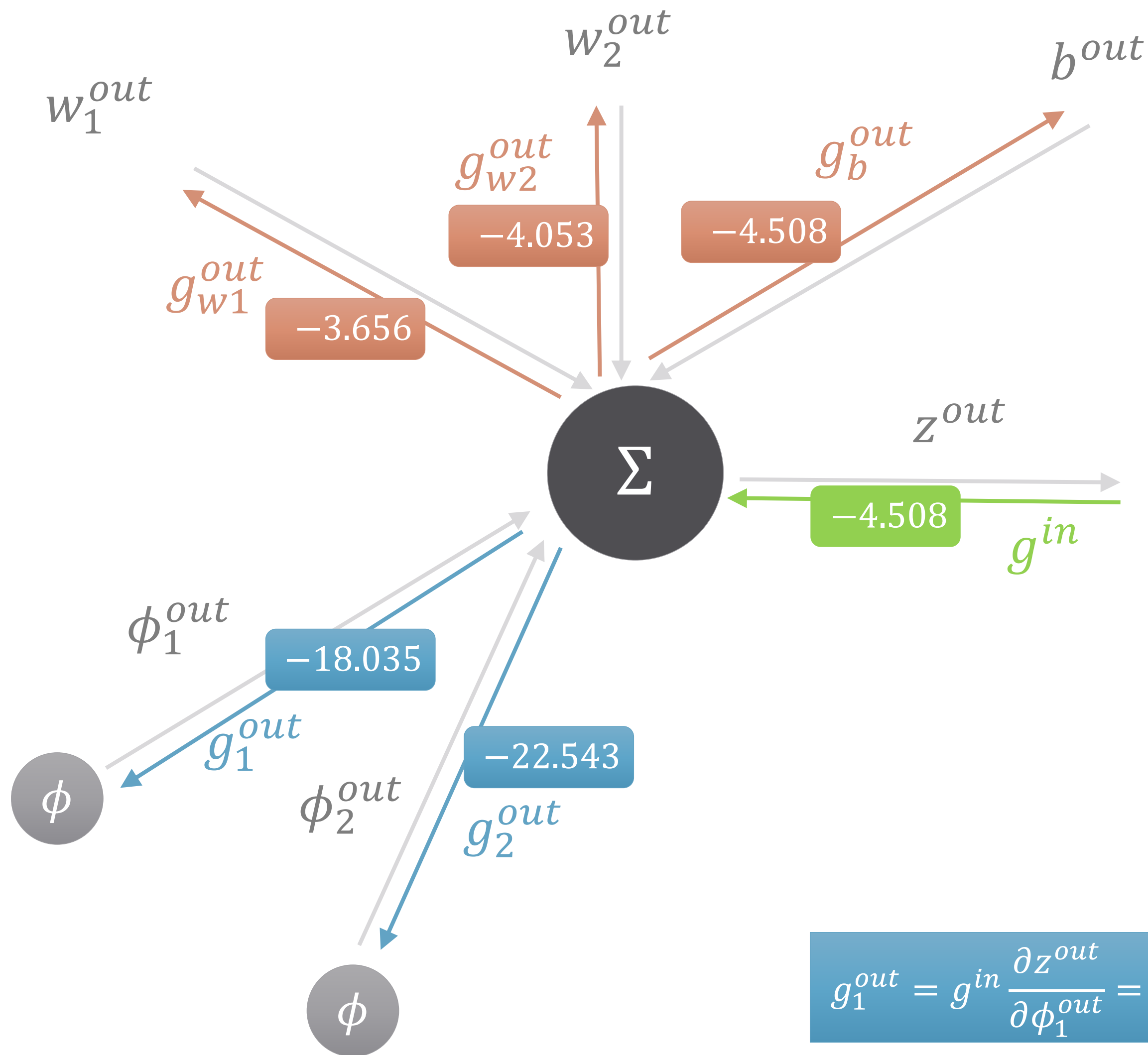$$\frac{\partial z^{out}}{\partial \phi_1^{out}} = w_1^{out} = 4$$

$$\frac{\partial z^{out}}{\partial w_1^{out}} = \phi_1^{out} = 0.811$$

$$\frac{\partial z^{out}}{\partial \phi_2^{out}} = w_2^{out} = 5$$

$$\frac{\partial z^{out}}{\partial w_2^{out}} = \phi_2^{out} = 0.899$$

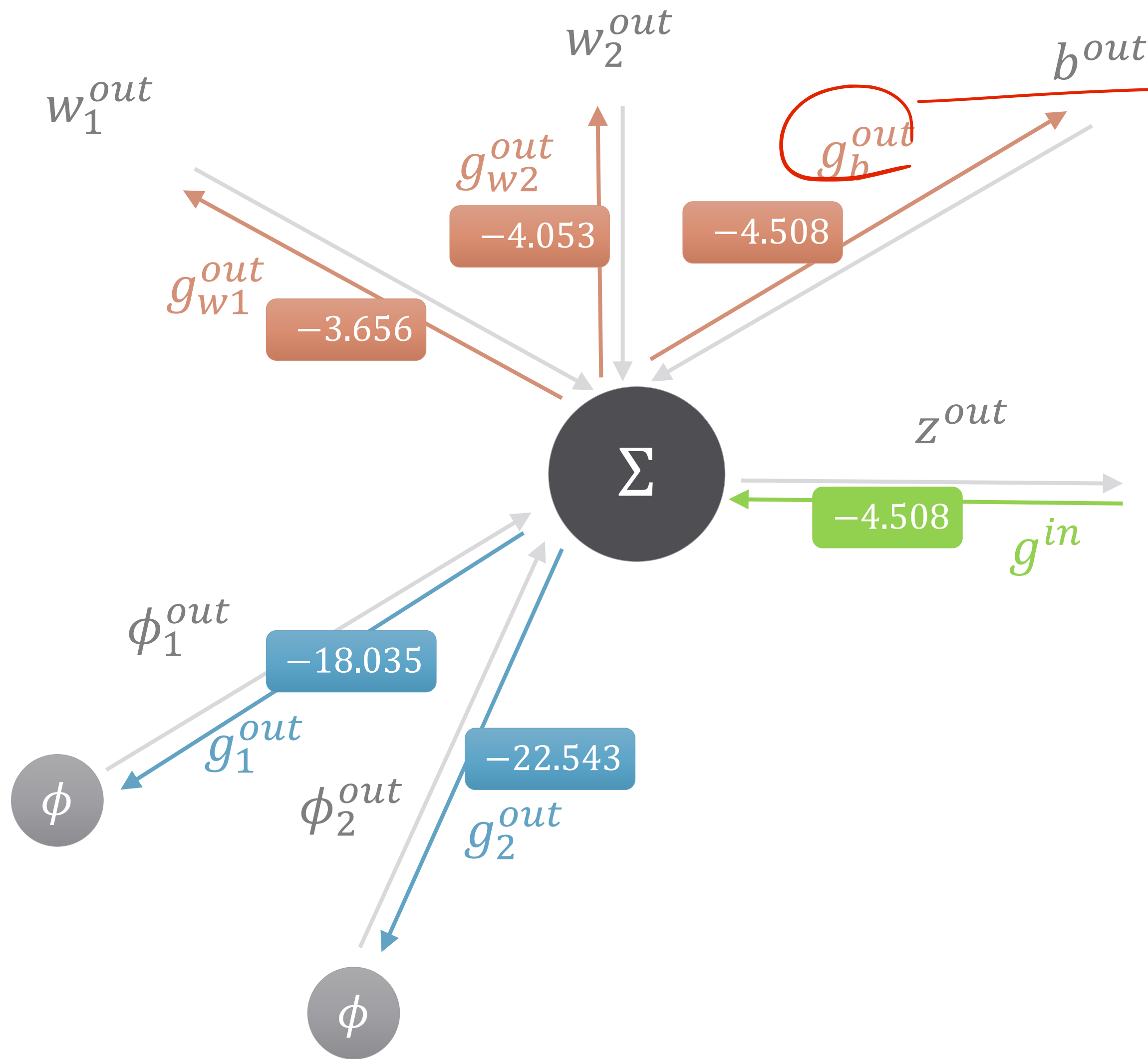$$\frac{\partial z^{out}}{\partial b^{out}} = 1$$

$$g_1^{out} = g^{in}\frac{\partial z^{out}}{\partial \phi_1^{out}} = -4.508 \times 4 = -18.035$$

$$g_{w1}^{out} = g^{in}\frac{\partial z^{out}}{\partial w_1^{out}} = -4.508 \times 0.811 = -3.656$$

$$g_2^{out} = g^{in}\frac{\partial z^{out}}{\partial \phi_2^{out}} = -4.508 \times 5 = -22.543$$

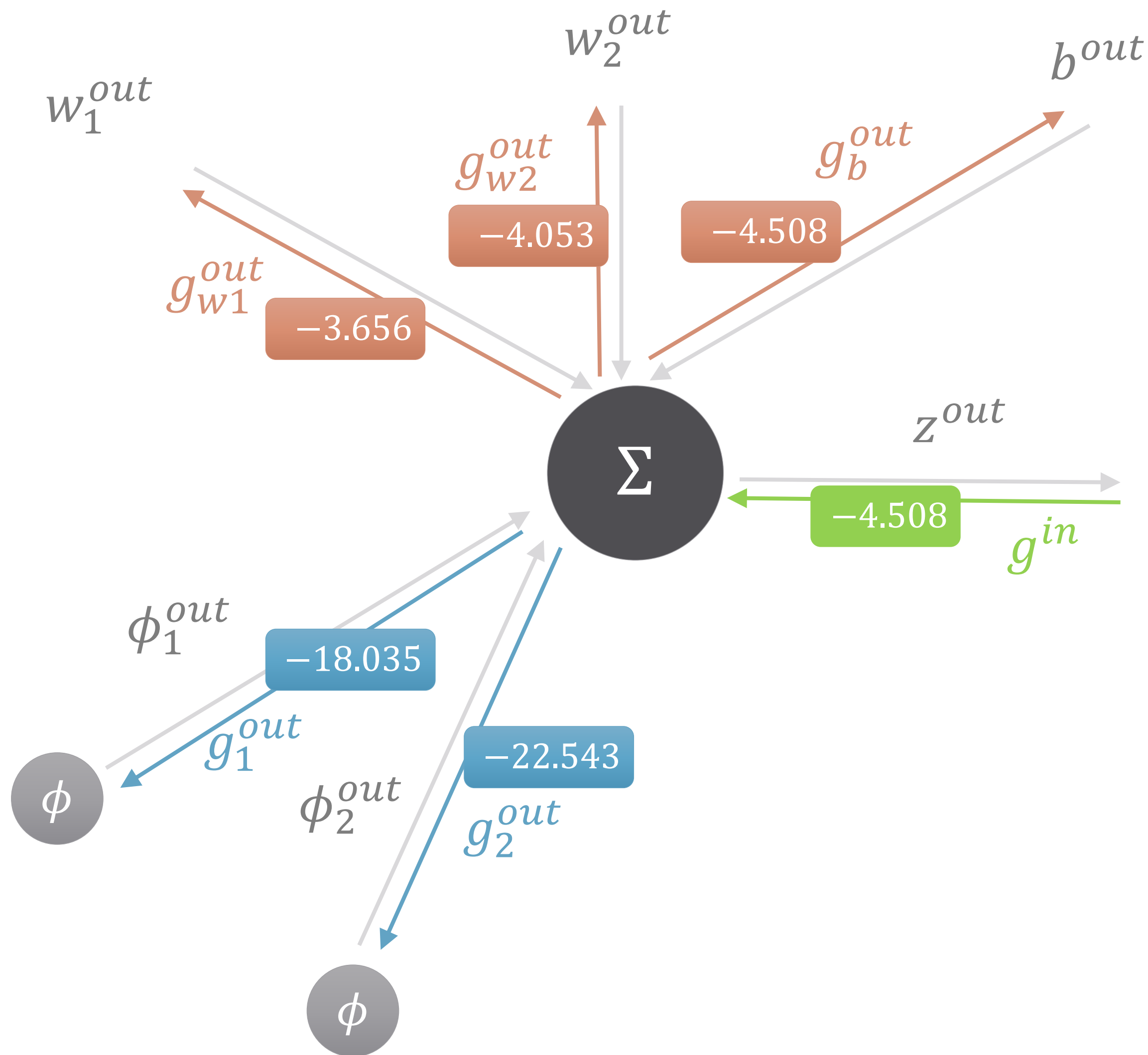$$g_{w2}^{out} = g^{in}\frac{\partial z^{out}}{\partial w_2^{out}} = -4.508 \times 0.899 = -4.053$$

$$g_b^{out} = g^{in}\frac{\partial z^{out}}{\partial b^{out}} = -4.508 \times 1 = -4.508$$

$w_1^{out}$

$w_2^{out}$

$b^{out}$

$g_{w2}^{out}$

$g_b^{out}$

−4.053

−4.508

$g_{w1}^{out}$

−3.656

Σ

$z^{out}$

−4.508   $g^{in}$

$\phi_1^{out}$

−18.035

$g_1^{out}$

$\phi$

$\phi_2^{out}$

−22.543

$g_2^{out}$

$\phi$

*Update Parameters*
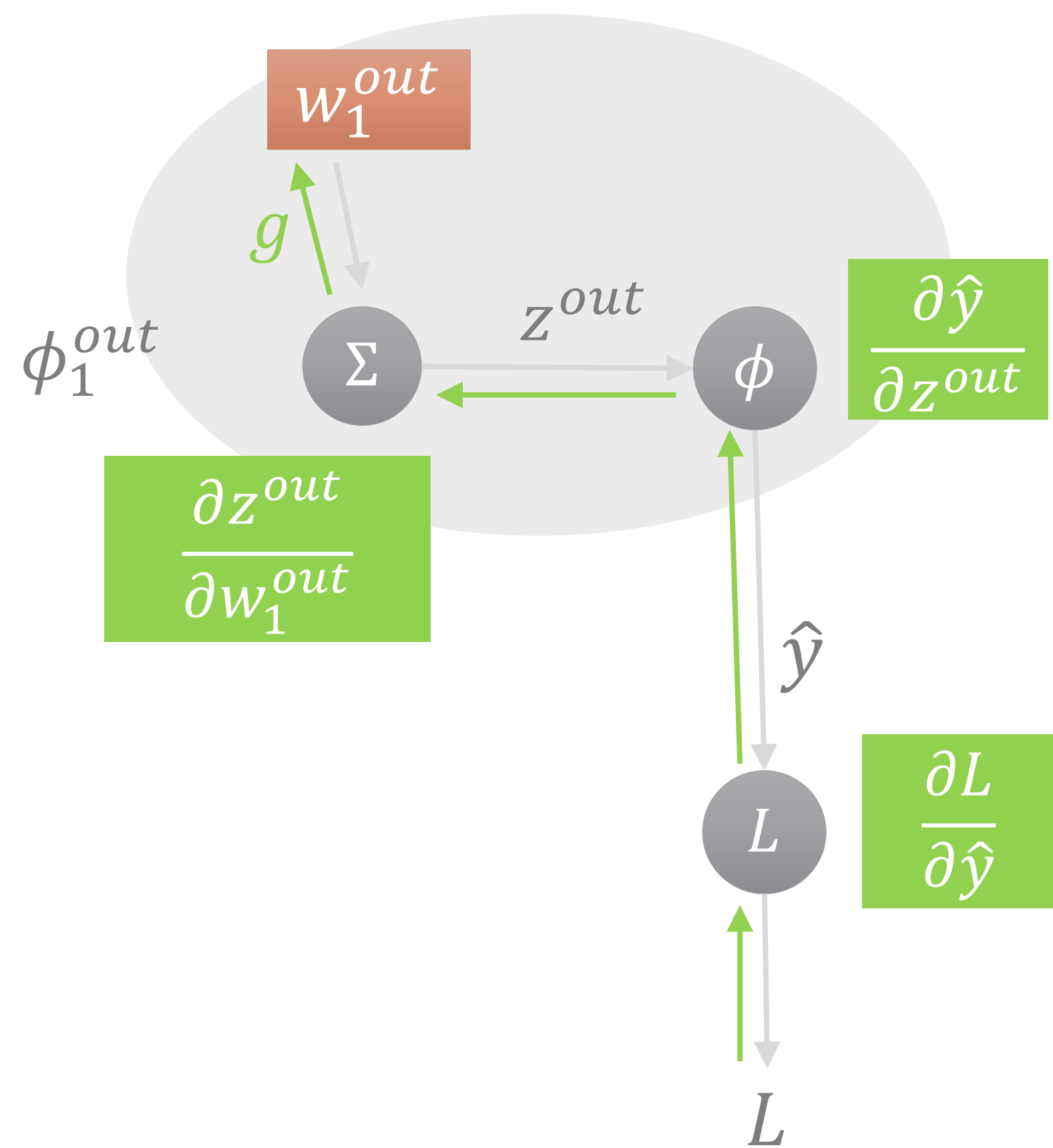
$$w^{new} = w^{old} - (\eta \times g)$$

Learning Rate
(0.1)

Update Parameters

$$w^{new} = w^{old} - (\eta \times g)$$

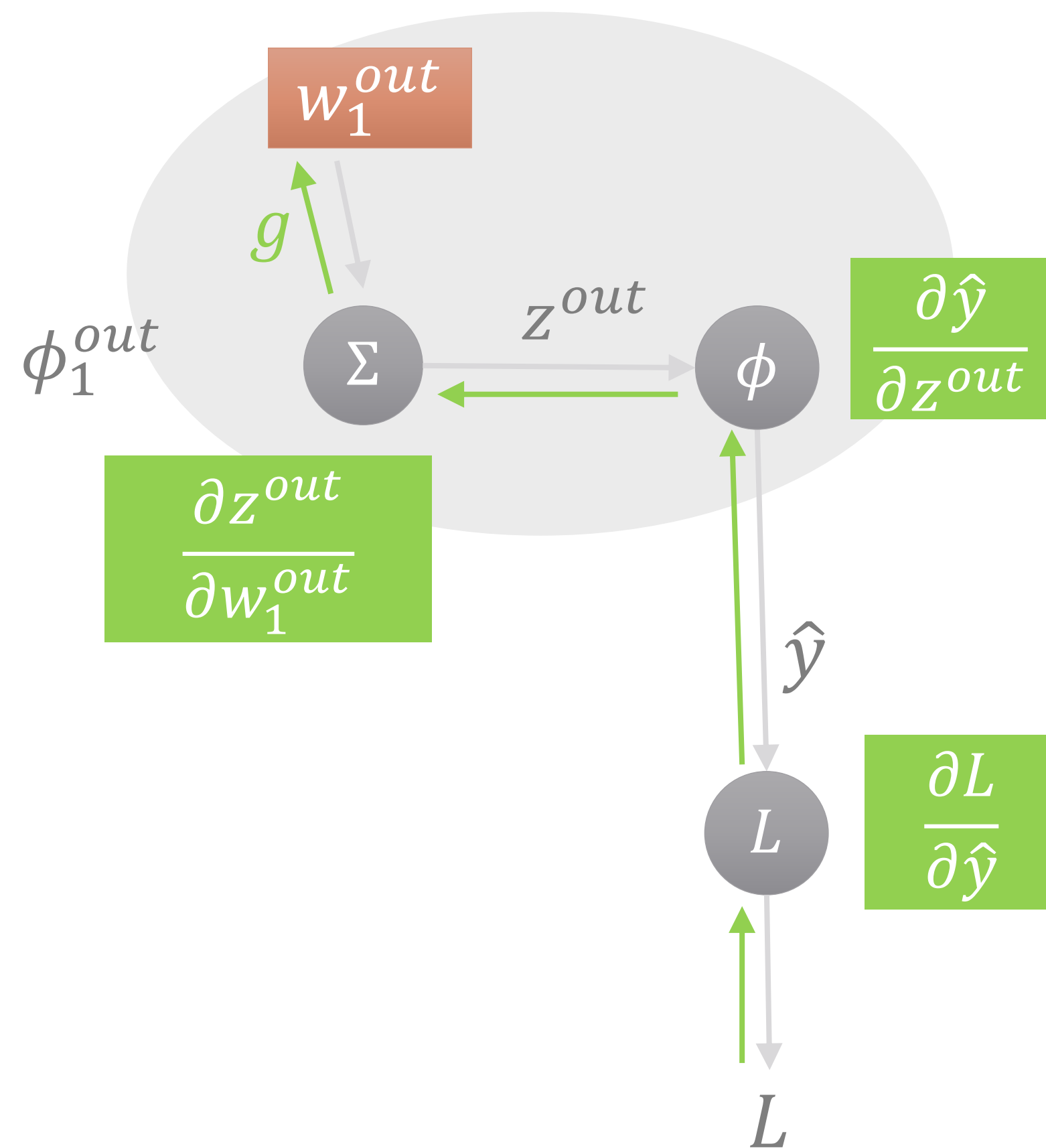Learning Rate
(0.1)

| Variable | Init Val | $\eta \times g$ | Updated Val |
|----------|----------|-----------------|-------------|
| $w_1^{out}$ | 4 | $-0.366$ | 4.366 |
| $w_2^{out}$ | 5 | $-0.405$ | 5.405 |
| $b^{out}$ | 0 | $-0.451$ | 0.451 |

Chain Rule

$$g = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{out}} \frac{\partial z^{out}}{\partial w_1^{out}} = \frac{\partial L}{\partial w_1^{out}}$$

## Chain Rule

$$g = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{out}} \frac{\partial z^{out}}{\partial w_1^{out}} = \frac{\partial L}{\partial w_1^{out}}$$
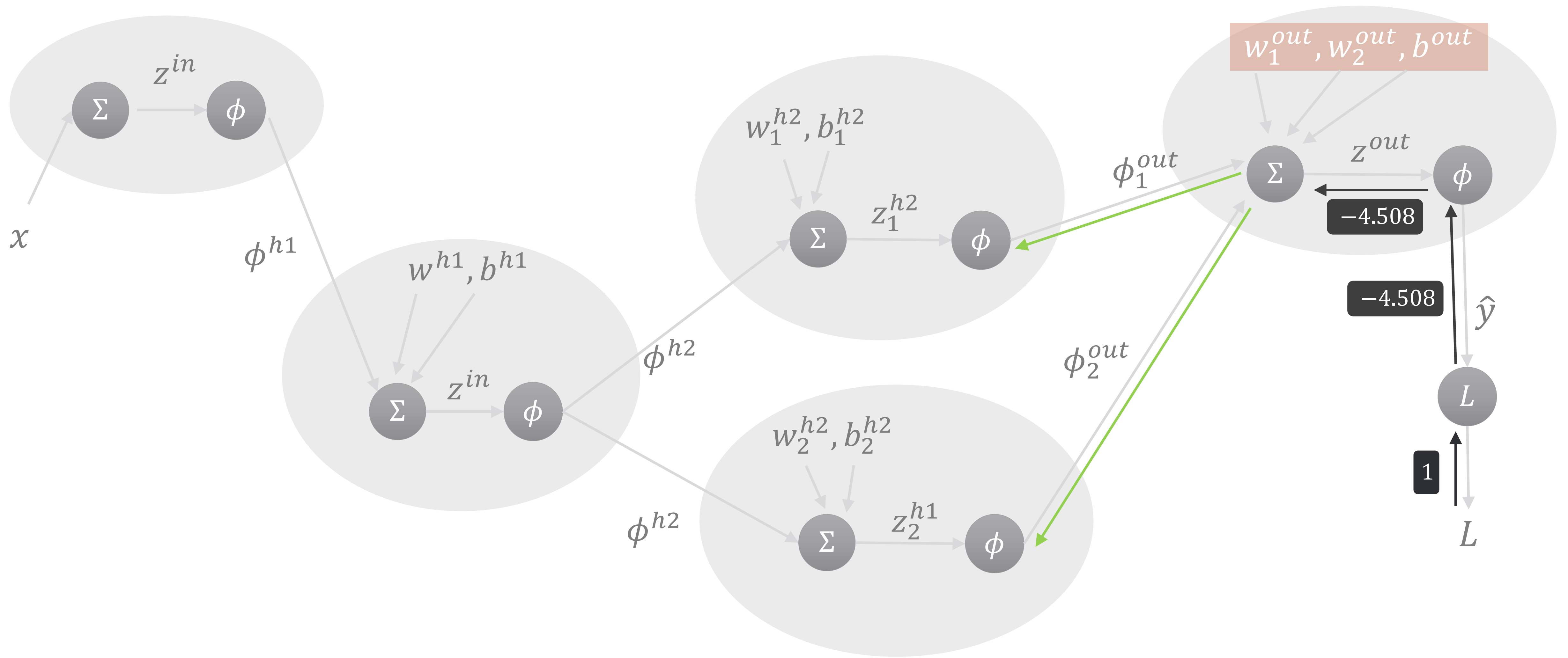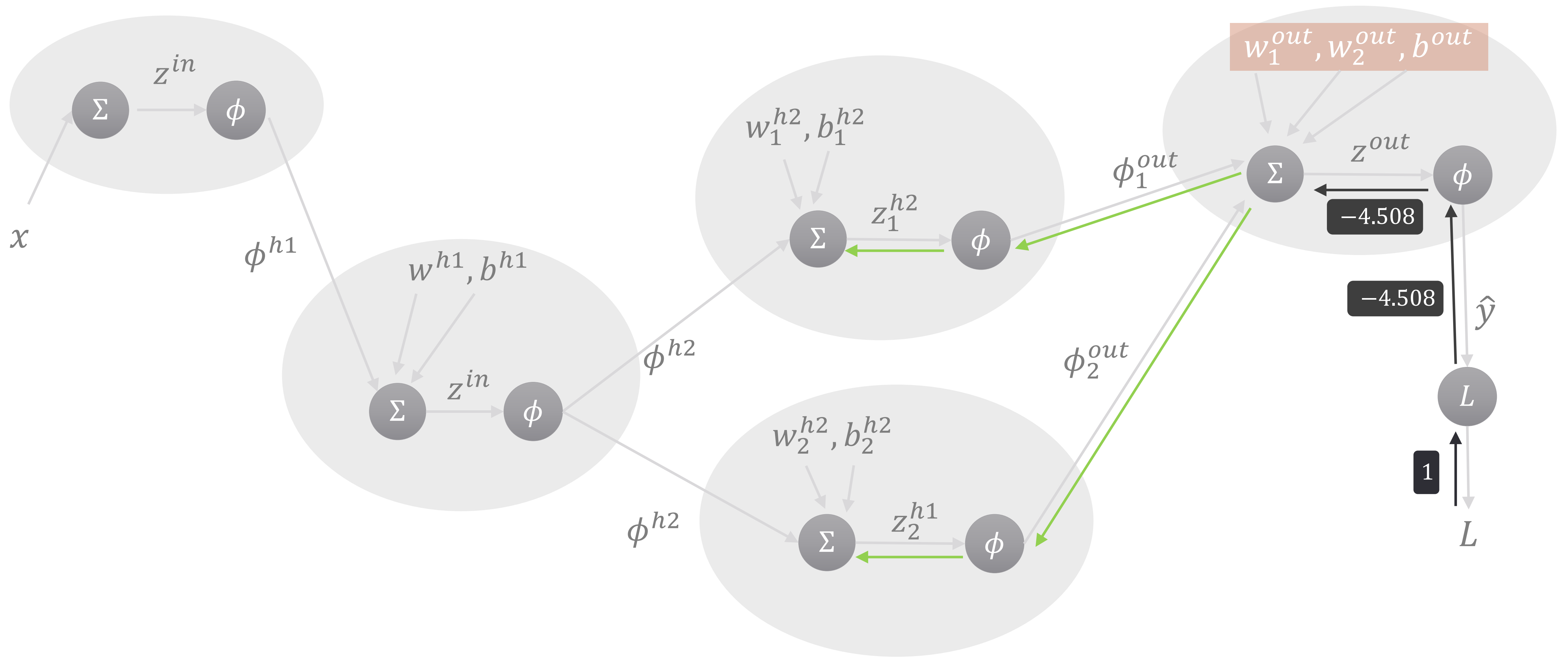
## Previously

$$w^{new} = w^{old} - (\eta \times g)$$

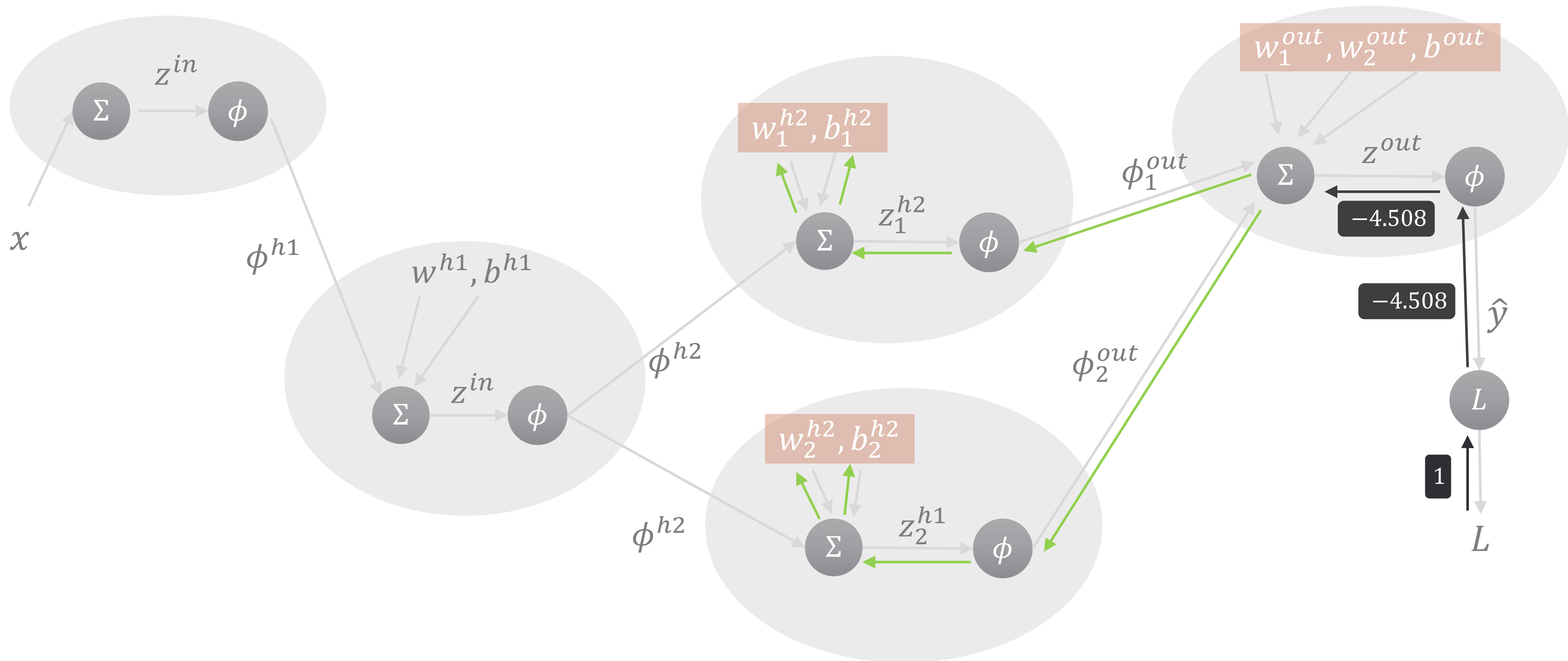## Formally

$$w^{new} = w^{old} - \eta \frac{\partial L}{\partial w_1^{out}}$$

## Gradient Descent

$z^{in}$

$x$

$\phi^{h1}$

$w^{h1}, b^{h1}$

$z^{in}$

$\phi^{h2}$

$w_1^{h2}, b_1^{h2}$

$z_1^{h2}$

$\phi^{h2}$

$w_2^{h2}, b_2^{h2}$

$z_2^{h1}$

$w_1^{out}, w_2^{out}, b^{out}$

$\phi_1^{out}$

$z^{out}$

$-4.508$

$-4.508$

$\phi_2^{out}$

$\hat{y}$

$L$

$1$

$L$

# Updated Values

| Variable | Init Val | Updated Val |
| --- | --- | --- |
| $w^{h1}$ | 1 | 1.122 |
| $b^{h1}$ | 0 | 0.228 |
| $w_1^{h2}$ | 2 | 2.201 |
| $b_1^{h2}$ | 0 | 0.275 |
| $w_2^{h2}$ | 3 | 3.148 |
| $b_2^{h2}$ | 0 | 0.203 |
| $w_1^{out}$ | 4 | 4.366 |
| $w_2^{out}$ | 5 | 5.405 |
| $b^{out}$ | 0 | 0.451 |

# Updated Prediction/Loss

Observed: $x = 1, y = 10$

| Variable | Before | After |
|:---:|:---:|:---:|
| $\hat{y}$ | 7.745 | 9.798 |
| $L$ | 5.082 | 0.352 |

# Multiple Inputs and Outputs

$x_1 \longrightarrow$ In 1

$x_2 \longrightarrow$ In 2

$x_3 \longrightarrow$ In 3

Out 1 $\longrightarrow \hat{y}_1$

Out 2 $\longrightarrow \hat{y}_2$

Out 3 $\longrightarrow \hat{y}_3$

# Activation Functions

**Tanh**

ออก!

เข้า !

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

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$$

V7 Labs

https://www.v7labs.com/blog/neural-networks-activation-functions

## ReLU



$$f(x) = max\,(0, x)$$

## Leaky ReLU

max(0.1 * x,x)

max(0.1 * x,x)



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

# Optimizer

- Stochastic gradient descent (SGD)
  - Estimate the actual gradient (calculated from the entire data set) by a value from subset of data (batch).

- Adaptive moment estimation (ADAM)
  - Stochastic gradient descent with adaptive learning rate optimization algorithm

# Classification

- Output nodes equal to number of class.

  - One-hot encoding

- Use *softmax* function to calculate probability of each class.

  - $p_j = softmax(\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_C) = \dfrac{e^{\hat{y}_j}}{\Sigma_{k=1}^{C} e^{\hat{y}_k}}$

- Loss function

  - *Categorical cross entropy (CCE)*

  - $CCE = -\dfrac{1}{N} \Sigma_{i=1}^{N} \Sigma_{k=1}^{C} \mathcal{X}_{y_k \in C_k} \ln(p_k)$

# CCE Example
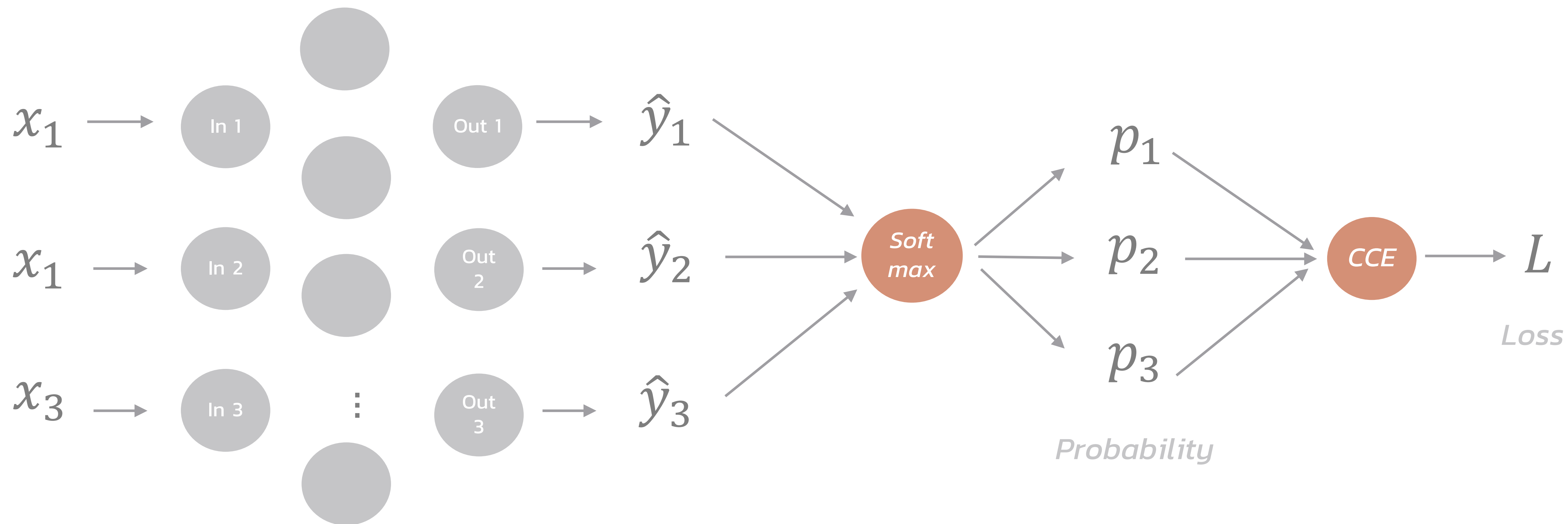
$$y_1 = [0, 1, 0] \qquad\qquad y_2 = [0, 0, 1]$$

$$p_1 = [0.05, 0.95, 0] \qquad p_2 = [0.1, 0.8, 0.1]$$

$$CCE = -\frac{1}{2}(\ln 0.95 + \ln 0.1) = 1.177$$

$x_1 \rightarrow$ In 1

$x_1 \rightarrow$ In 2

$x_3 \rightarrow$ In 3

Out 1 $\rightarrow \hat{y}_1$

Out 2 $\rightarrow \hat{y}_2$

Out 3 $\rightarrow \hat{y}_3$

Soft max

$p_1$

$p_2$

$p_3$

Probability

CCE $\rightarrow L$

Loss

# Loss function

| Loss function | Usage | Examples | |
|---|---|---|---|
| | | **Using probabilities** *from_logits=False* | **Using logits** *from_logits=True* |
| BinaryCrossentropy | Binary classification | y_true: 1 <br> y_pred: 0.69 | y_true: 1 <br> y_pred: 0.8 |
| CategoricalCrossentropy | Multiclass classification | y_true: 0  0  1 <br> y_pred: 0.30  0.15  0.55 | y_true: 0  0  1 <br> y_pred: 1.5  0.8  2.1 |
| Sparse CategoricalCrossentropy | Multiclass classification | y_true: 2 <br> y_pred: 0.30  0.15  0.55 | y_true: 2 <br> y_pred: 1.5  0.8  2.1 |

We will use this one

# Regression

- Output layer
  - No *"softmax"*
- Loss

  Mean squared error
  - Mean absolute (percentage) error
- *Scale both X and y data*
  - Scaling X: more stable model (small weights)
  - Scaling y: matching output of activation function / smaller gradient