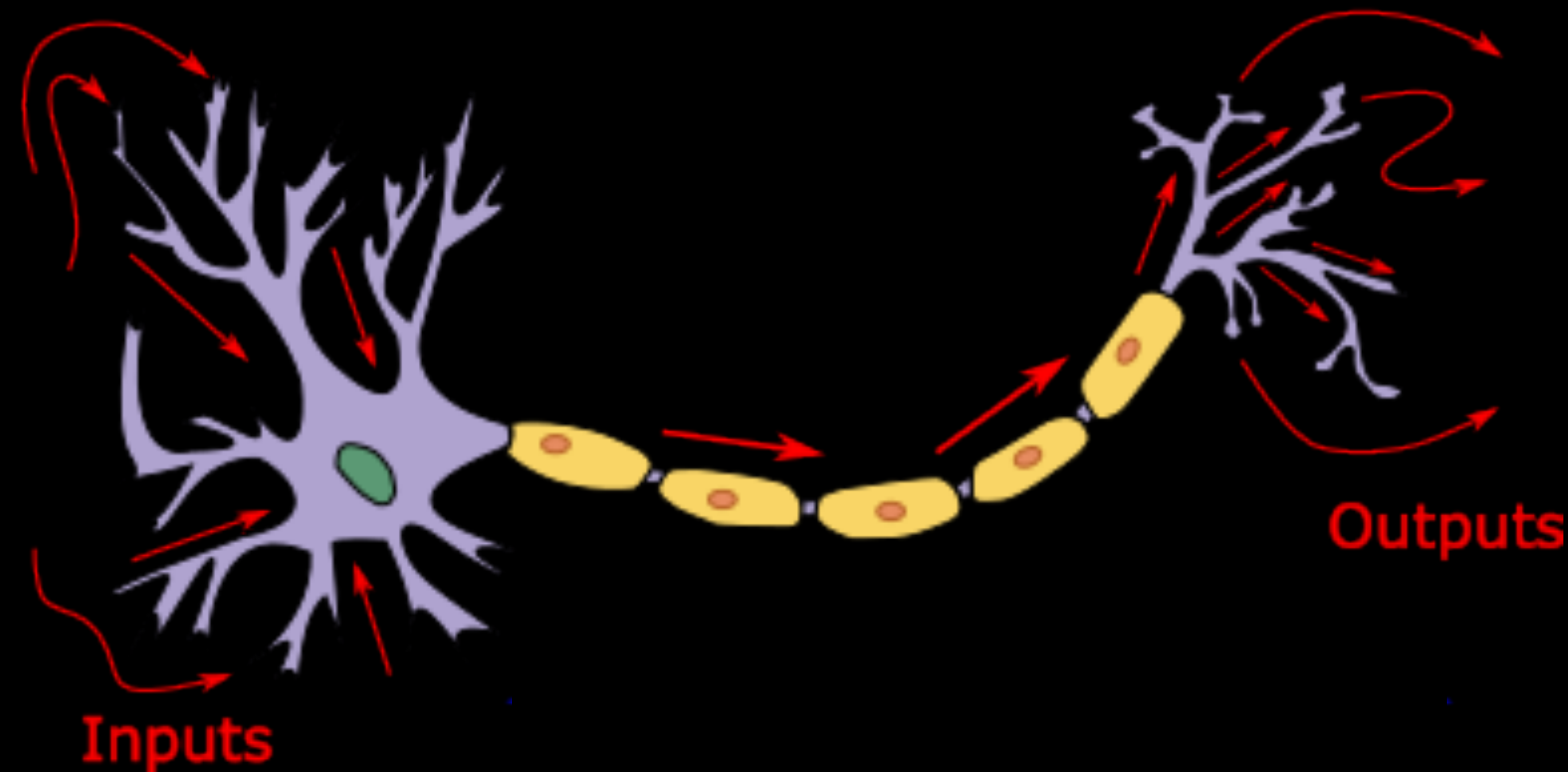
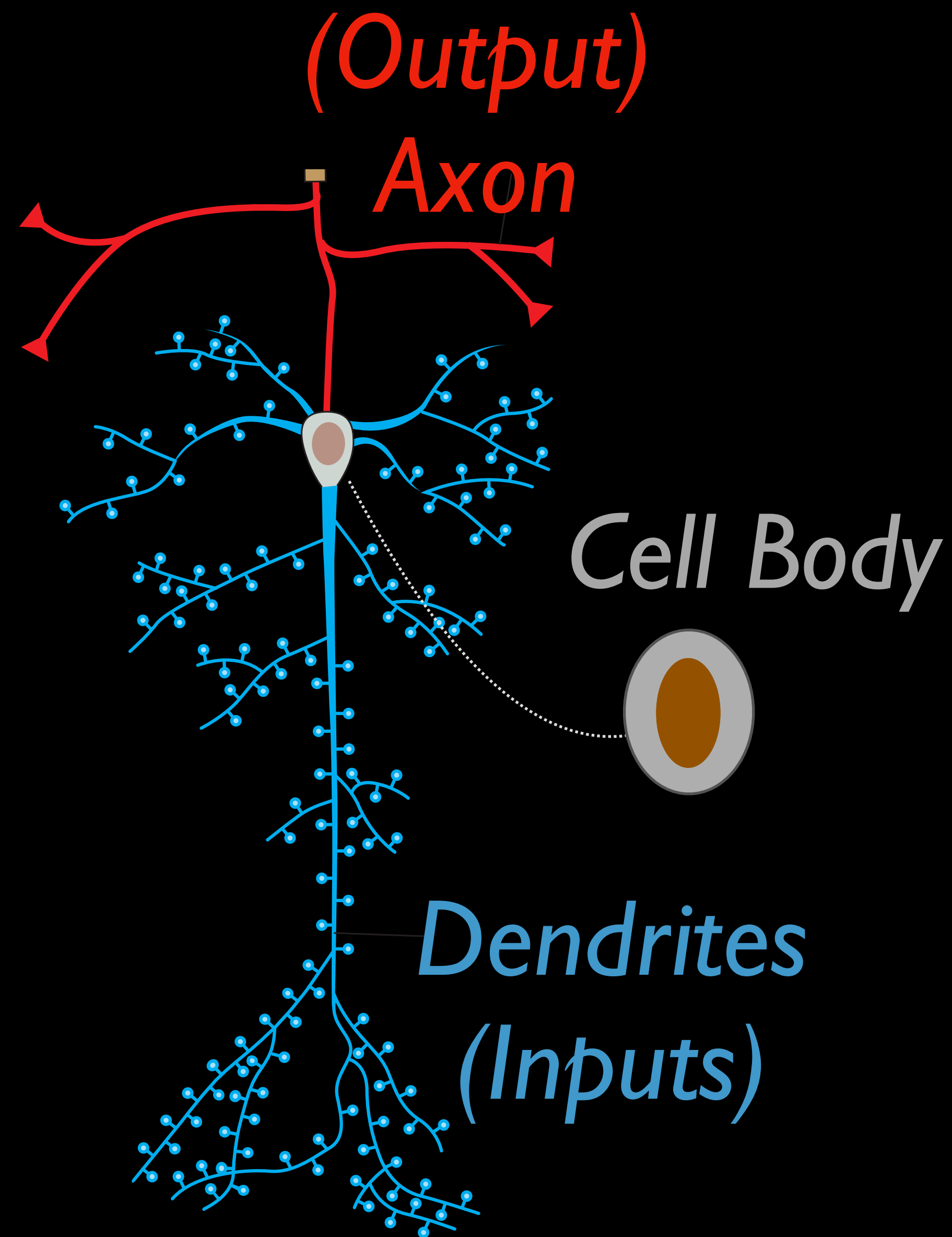


COS234: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer



Topic: Feed-Forward Neural Networks



“Neuron”

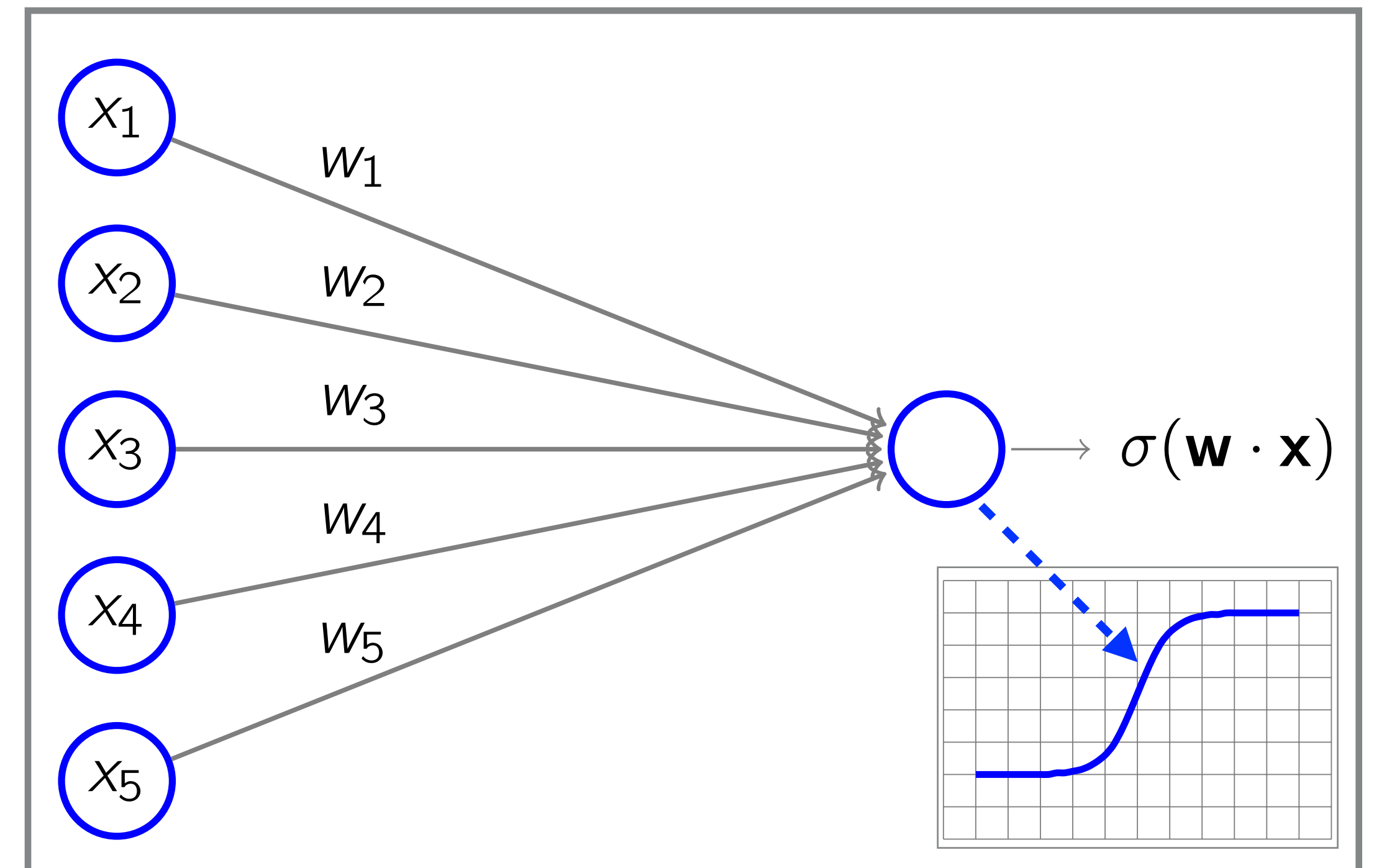
GLM \Rightarrow Neuron

Loss aside, GLM is input-output map:

$$\hat{y} = h(\mathbf{w} \cdot \mathbf{x})$$

Encapsulate as “black box” called Neuron

Reuse neurons with different parameters



Two Layer Neural Network

“Concatenate” multiple mappings:

$$h_1(\mathbf{w}_1 \cdot \mathbf{x}), h_2(\mathbf{w}_2 \cdot \mathbf{x}), \dots, h_m(\mathbf{w}_m \cdot \mathbf{x})$$

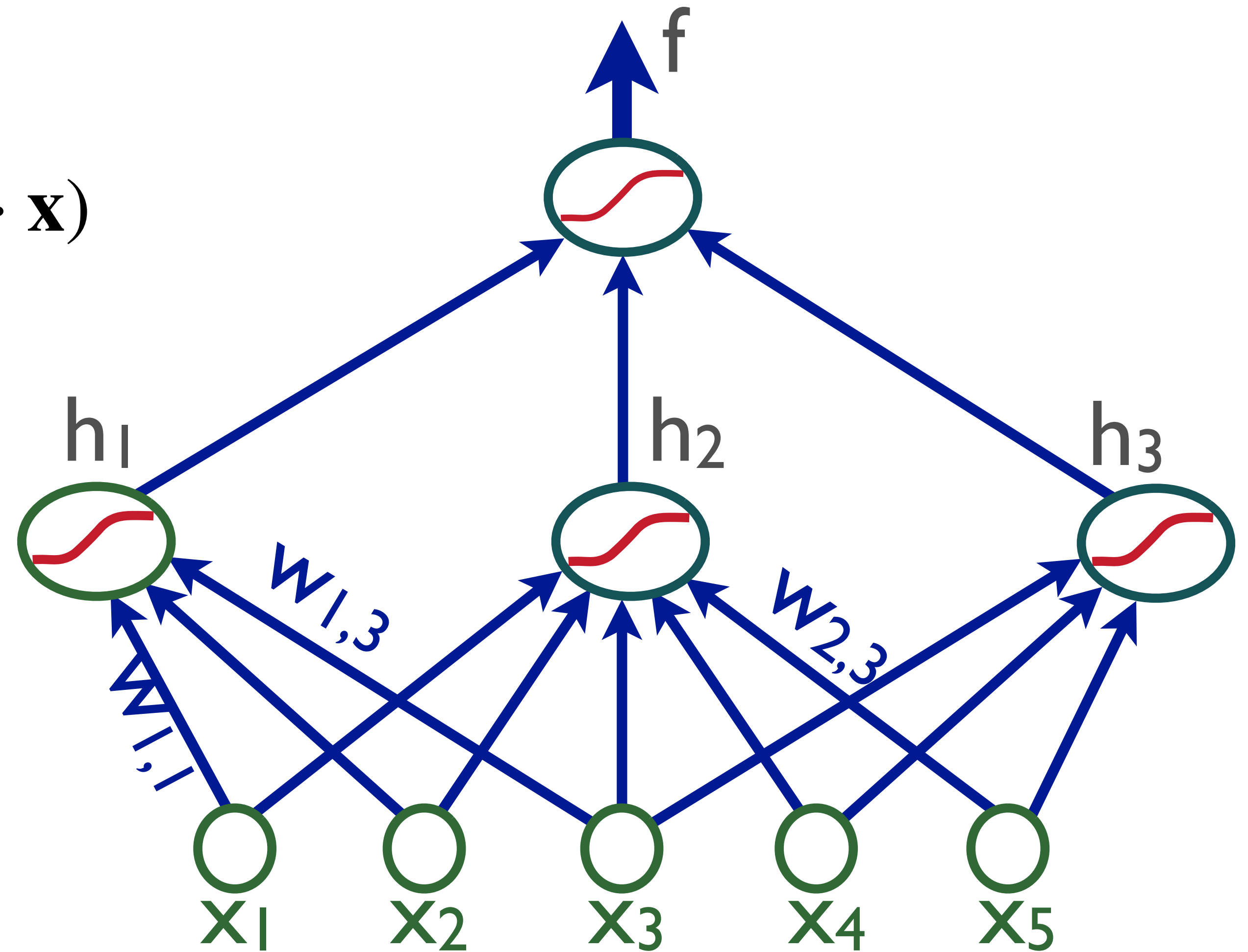
Activations need not be identical,
but typically: $h_j = h^{\text{generic}}$

Intermediate output is called a layer:

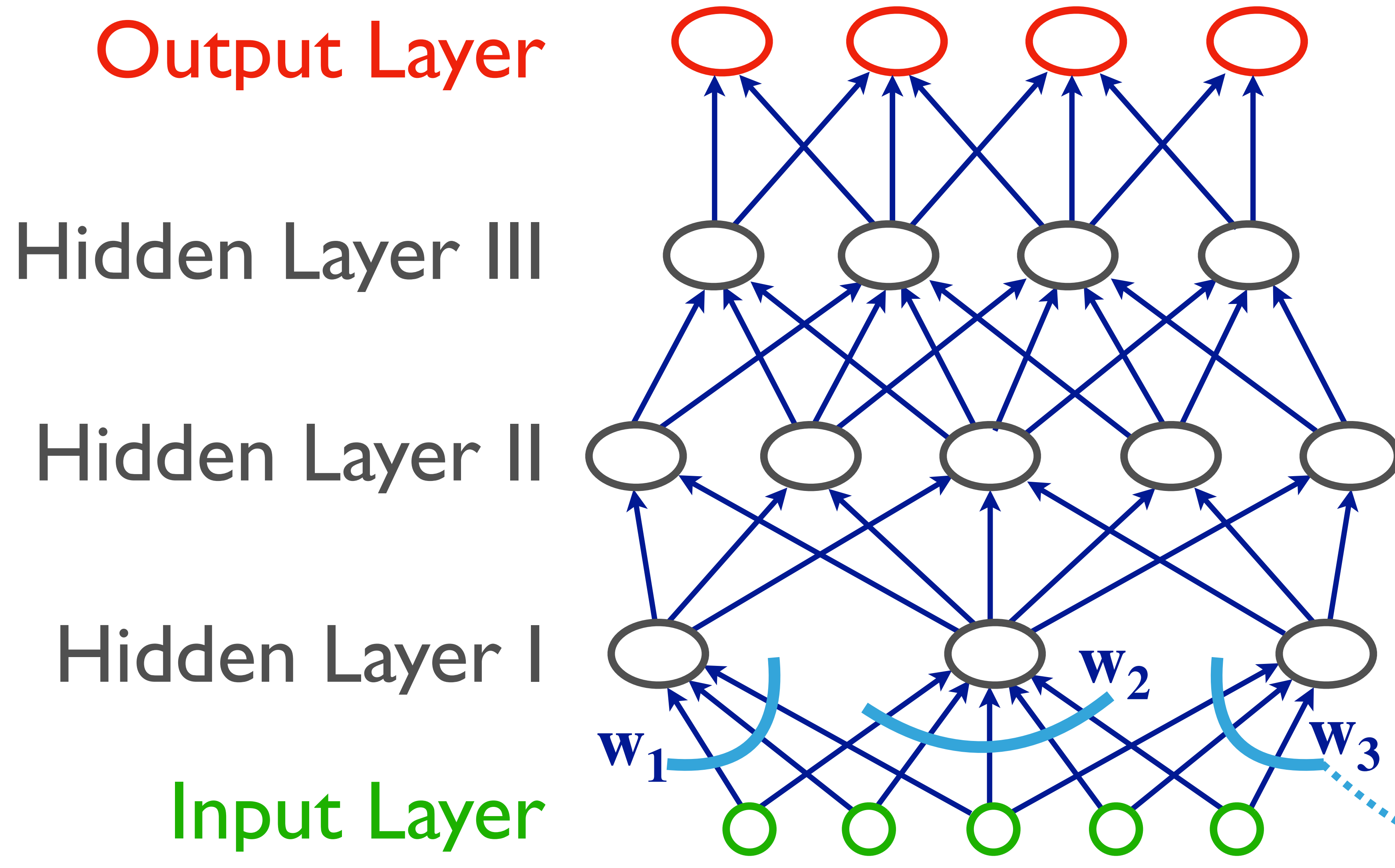
$$\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_m(\mathbf{x}))$$

Output of network:

$$\mathbf{f}(\mathbf{x}) = \mathbf{u} \cdot \mathbf{h}(\mathbf{x})$$



Multilayer Net



$$A^4 = \begin{bmatrix} -\mathbf{r}_1- \\ -\mathbf{r}_2- \\ -\mathbf{r}_3- \\ -\mathbf{r}_4- \end{bmatrix} \in \mathbf{R}^{4 \times 4}$$

$$A^3 = \begin{bmatrix} -\mathbf{v}_1- \\ -\mathbf{v}_2- \\ -\mathbf{v}_3- \\ -\mathbf{v}_4- \end{bmatrix} \in \mathbf{R}^{4 \times 5}$$

$$A^2 = \begin{bmatrix} -\mathbf{u}_1- \\ -\mathbf{u}_2- \\ \dots \\ -\mathbf{u}_5- \end{bmatrix} \in \mathbf{R}^{5 \times 3}$$

$$A^1 = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ -\mathbf{w}_3- \end{bmatrix} \in \mathbf{R}^{3 \times 5}$$

Multilayer Net

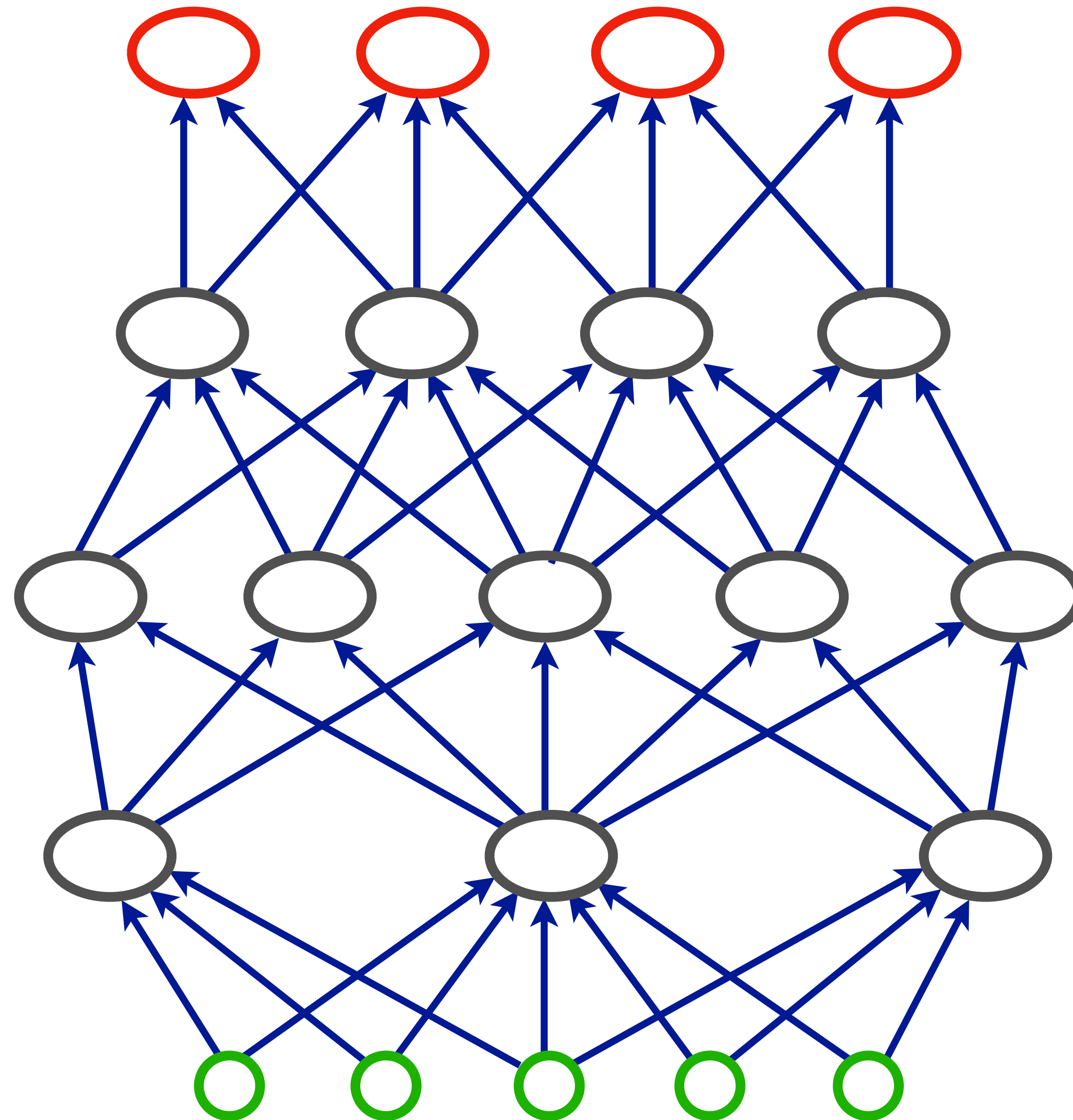
$$\mathbf{h}^4 \in \mathbf{R}^4$$

$$\mathbf{h}^3 \in \mathbf{R}^4$$

$$\mathbf{h}^2 \in \mathbf{R}^5$$

$$\mathbf{h}^1 \in \mathbf{R}^3$$

$$\mathbf{h}^0 \in \mathbf{R}^5$$



$$A^4 = \begin{bmatrix} -\mathbf{r}_1- \\ -\mathbf{r}_2- \\ -\mathbf{r}_3- \\ -\mathbf{r}_4- \end{bmatrix} \in \mathbf{R}^{4 \times 4}$$

$$A^3 = \begin{bmatrix} -\mathbf{v}_1- \\ -\mathbf{v}_2- \\ -\mathbf{v}_3- \\ -\mathbf{v}_4- \end{bmatrix} \in \mathbf{R}^{4 \times 5}$$

$$A^2 = \begin{bmatrix} -\mathbf{u}_1- \\ -\mathbf{u}_2- \\ \dots \\ -\mathbf{u}_5- \end{bmatrix} \in \mathbf{R}^{5 \times 3}$$

$$A^1 = \begin{bmatrix} -\mathbf{w}_1- \\ -\mathbf{w}_2- \\ -\mathbf{w}_3- \end{bmatrix} \in \mathbf{R}^{3 \times 5}$$

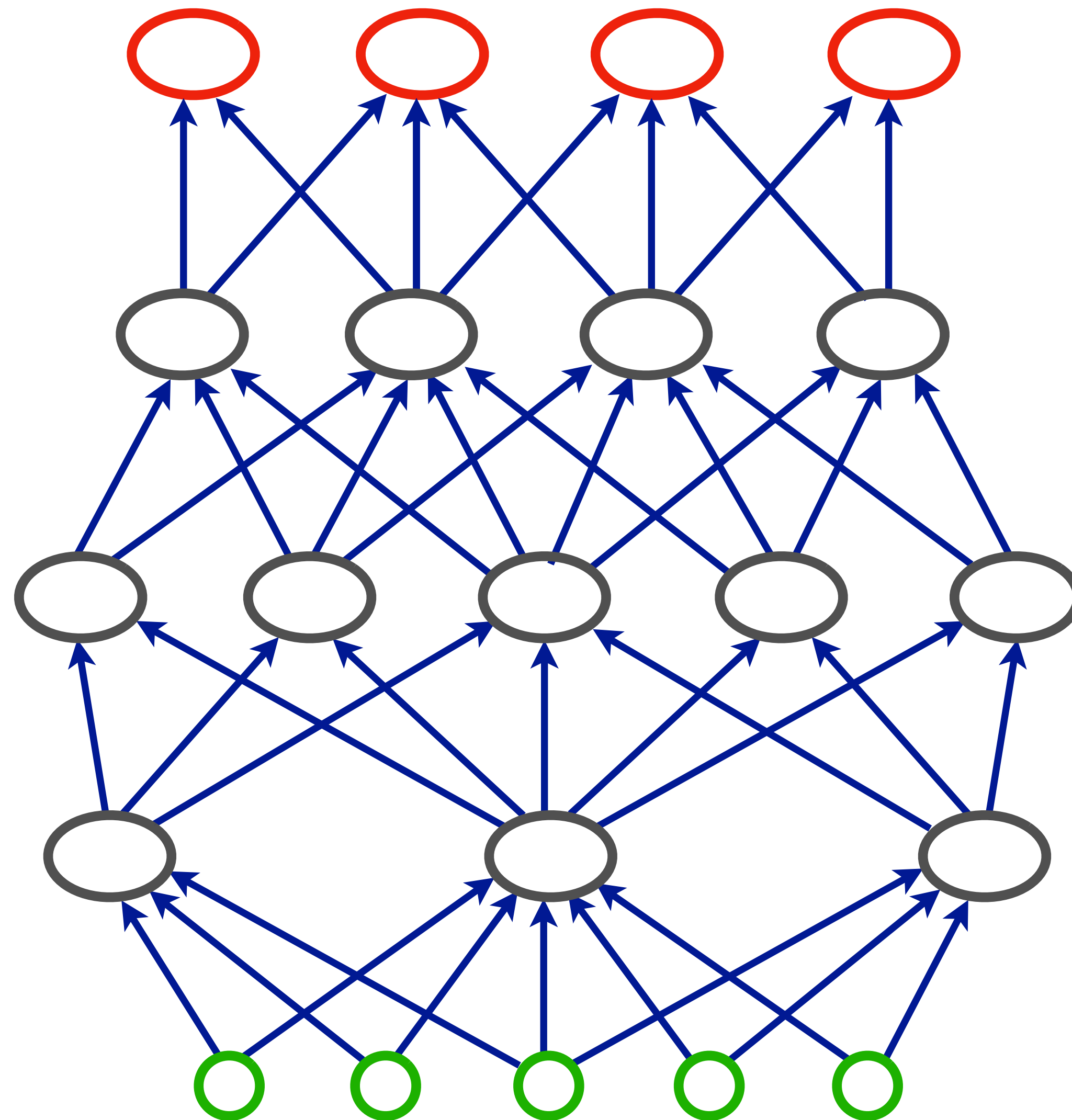
Multilayer Net

$$\mathbf{h}^4 = \sigma(A^4 \mathbf{h}^3)$$

$$\mathbf{h}^3 = \sigma(A^3 \mathbf{h}^2)$$

$$\mathbf{h}^2 = \sigma(A^2 \mathbf{h}^1)$$

$$\mathbf{h}^1 = \sigma(A^1 \mathbf{h}^0)$$

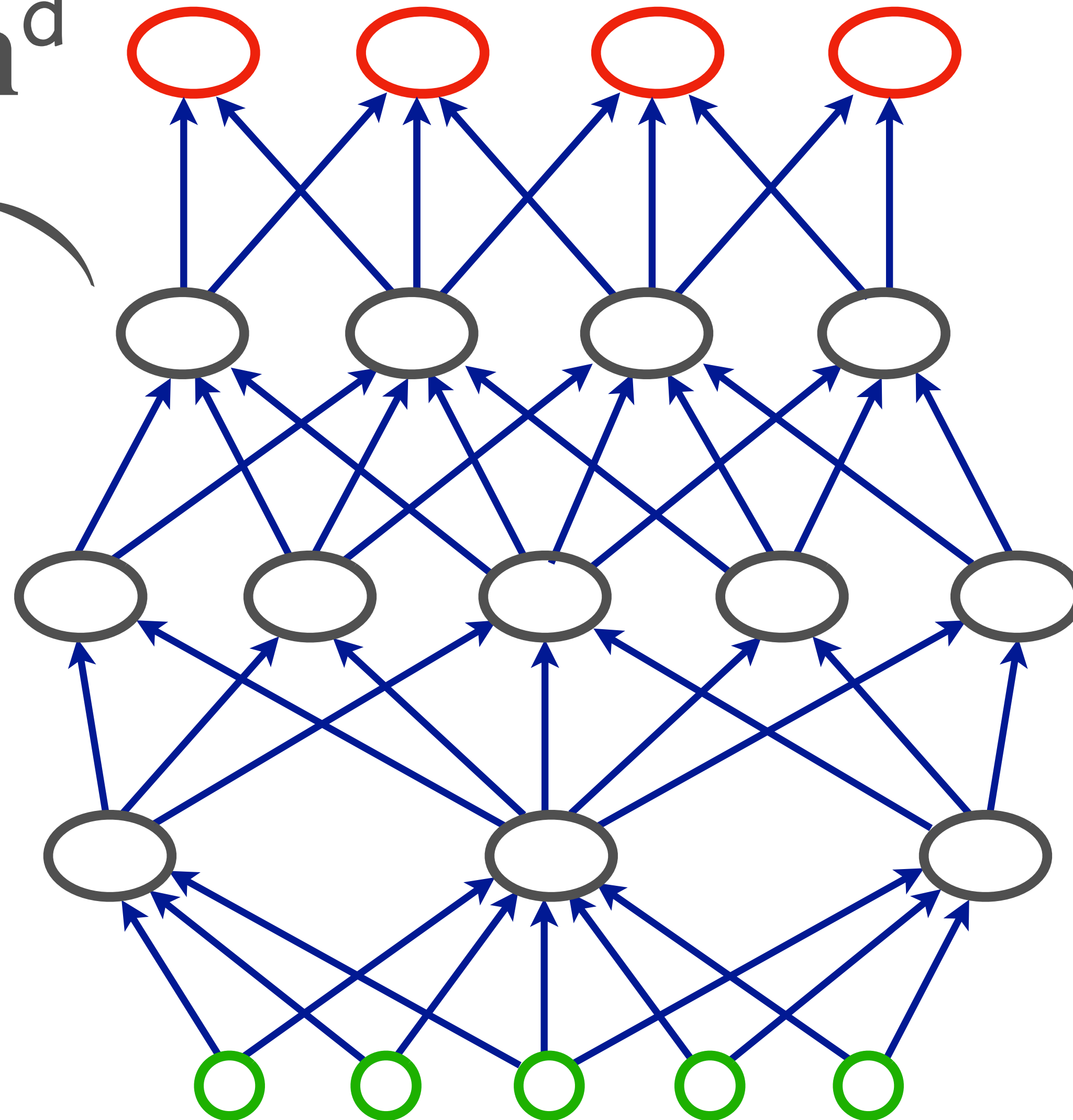


Inference aka Forward Pass

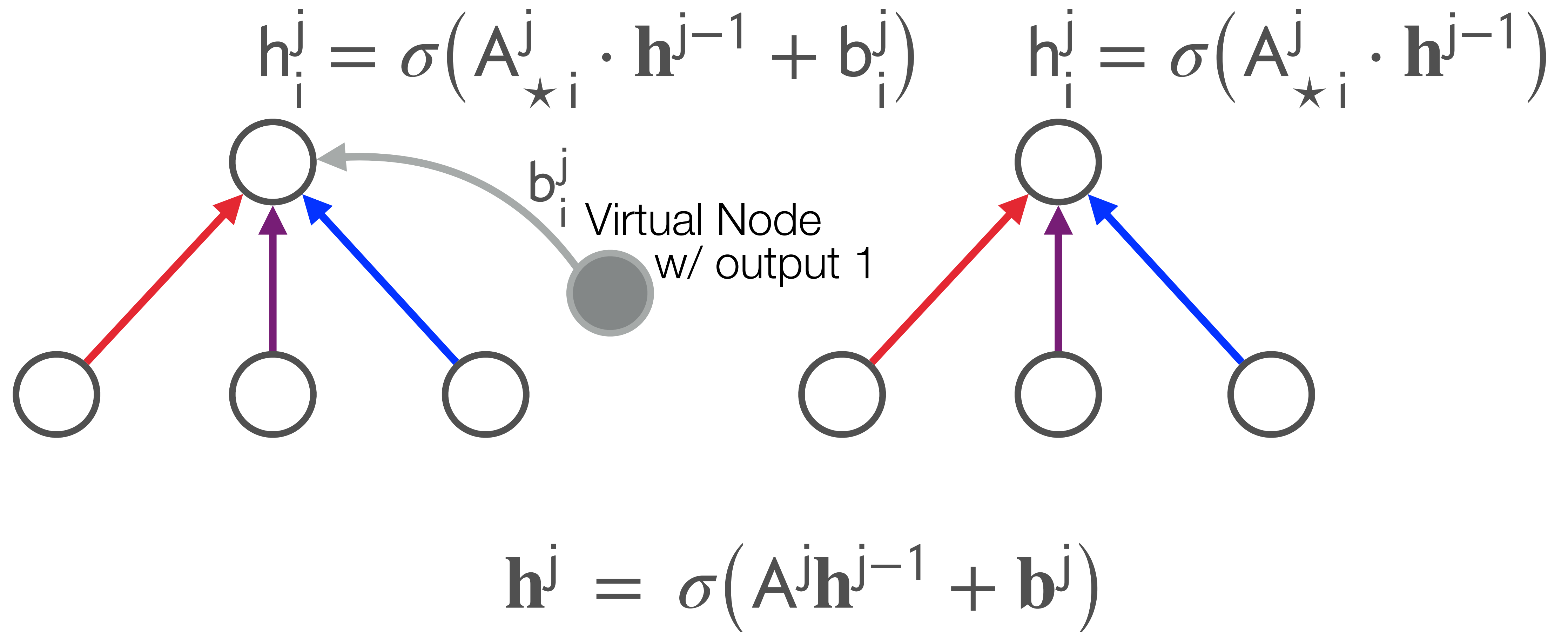
$$\mathbf{f}(\mathbf{x}) = \mathbf{h}^d$$

$$\mathbf{h}^i = \sigma(\mathbf{A}^i \mathbf{h}^{i-1})$$

$$\mathbf{h}^0 = \mathbf{x}$$

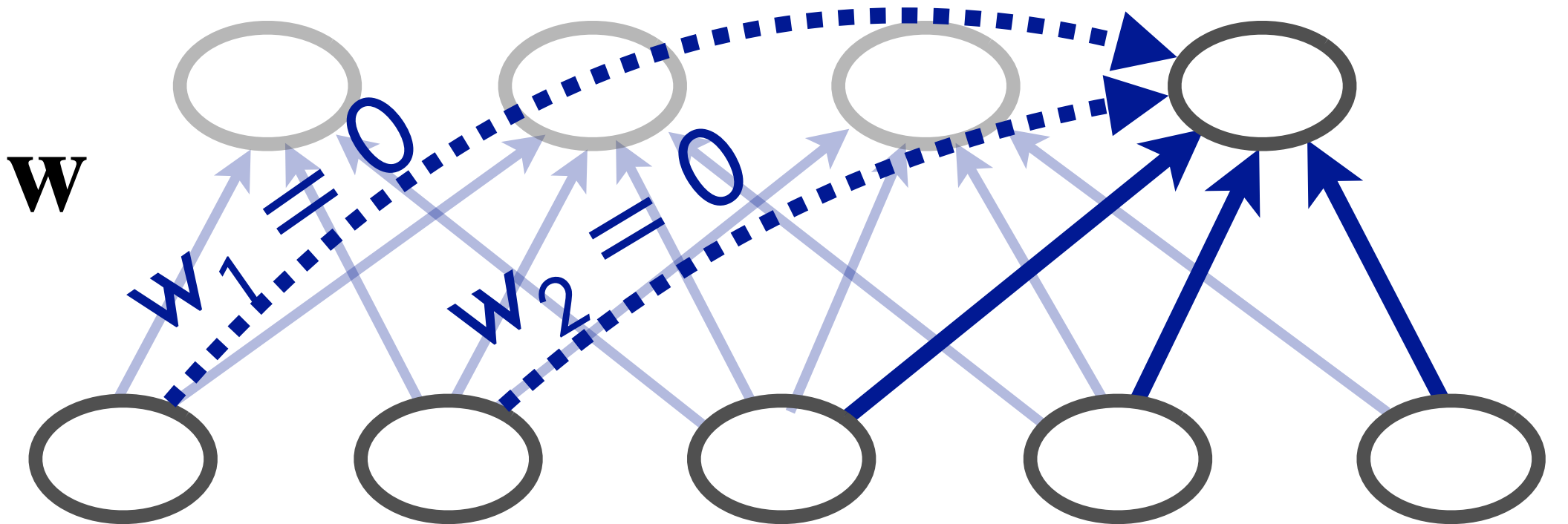


Adding Bias Vectors

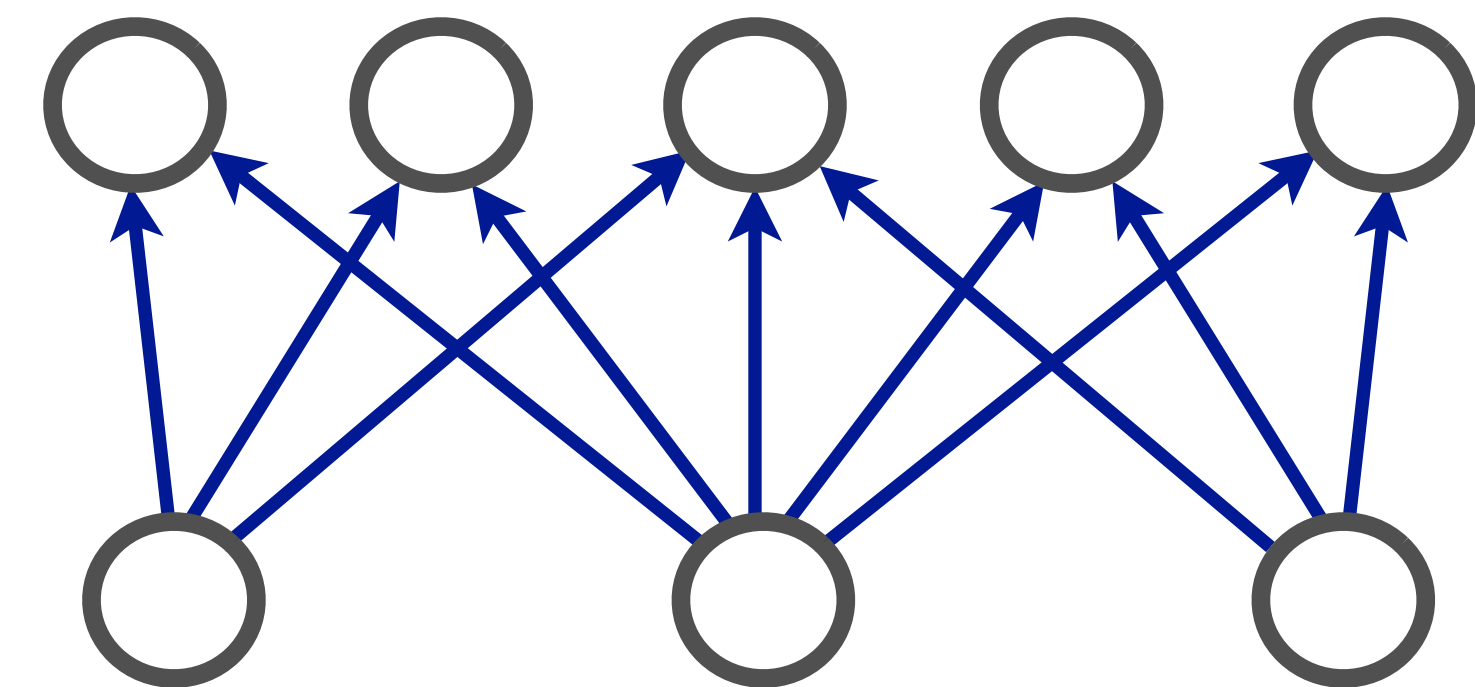
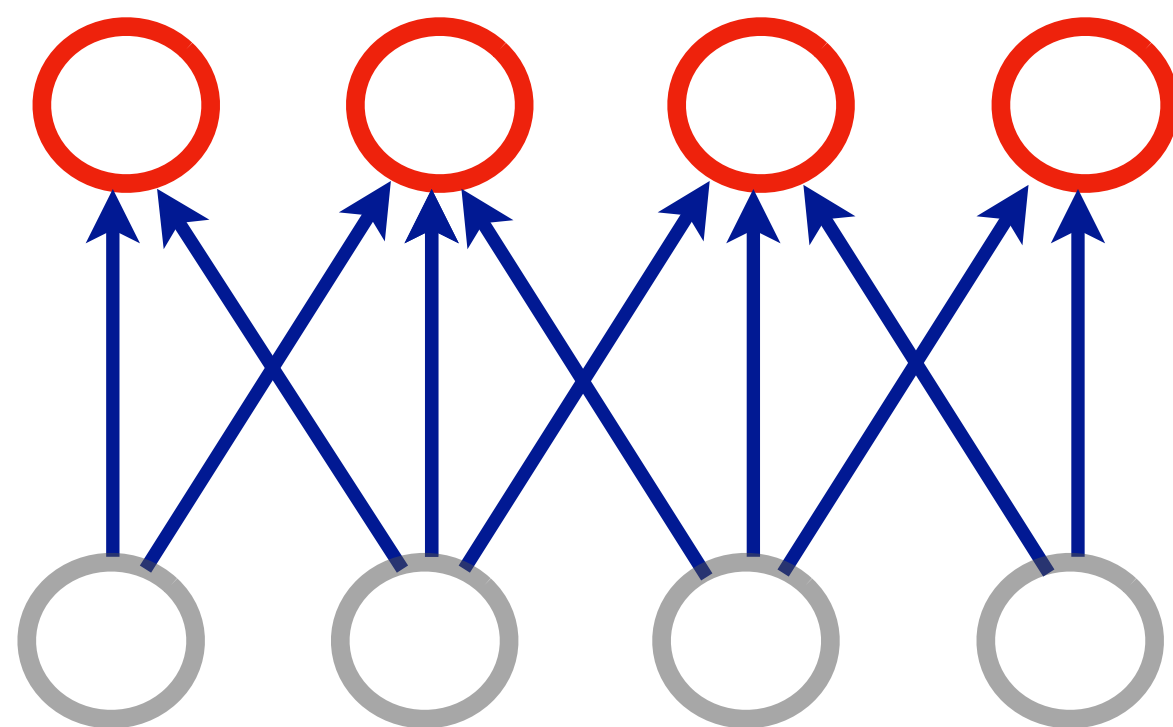


NN Architecture

Vector of weights of a single neuron \mathbf{w}
(incoming edges) can be sparse



Multiple Outputs $\mathbf{f}(\mathbf{x}) = \mathbf{h}^d$



#Neurons in each layer can be different

Feed-forward NN: Formal Definitions

- Obtained by connecting several neurons together
- **Feed-forward** networks: directed (acyclic) graph $G = (V, E)$ & denoting $m = |V|$
- Input nodes: nodes w/o incoming edges v_1, \dots, v_n
- Output nodes: nodes w/o outgoing edges v_{m-k+1}, \dots, v_m
- Weights associated with each edge $w : E \rightarrow \mathbb{R}$
- Computation defined by NN [for (input) node $j \in [n]$ we define $h[v_j] = x_j$]

$$h[v] = \sigma\left(\sum_{u \rightarrow v \in E} w[u \rightarrow v] h[u]\right)$$

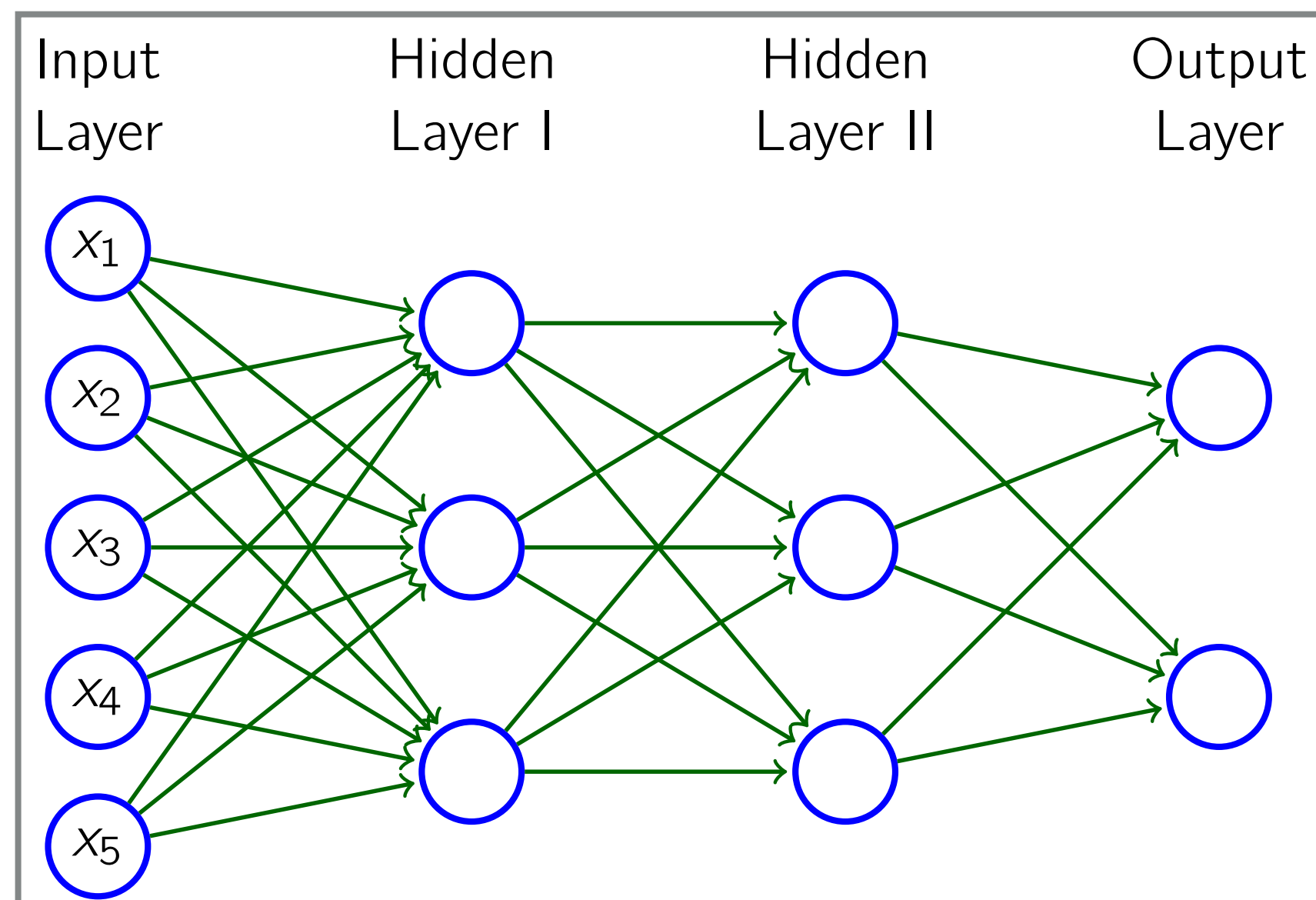
- Feed-forward NN defines a nonlinear mapping function $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$

Multilayer Perceptron (MLP)

Nodes divided into layers: $V = \left\{ V_j \right\}_{j=0}^d$

Edges solely between adjacent layers:

$$u \rightarrow v \in E \Rightarrow \text{for } j \in [d] : u \in V_{j-1} \wedge v \in V_j$$



#Inputs	$n=5$
Depth	$d=3$
#Neurons	$m = 5 + 3 + 3 + 2 = 13$
#Weights	$ E = 3 \times 5 + 3 \times 3 + 2 \times 3 = 30$

Multilayer Perceptron

Nodes divided into layers: $V = \left\{ V_j \right\}_{j=1}^d$

Edges solely between adjacent layers: $u \rightarrow v \in E \Rightarrow u \in V_j \wedge v \in V_{j+1}$

Define $A^j \in \mathbf{R}^{|V_{j-1}| \times |V_j|}$ where

$$A^j[\mathcal{J}(v_a), \mathcal{J}(v_b)] = w[v_a \rightarrow v_b] \text{ and } \mathcal{J}(v_a) = a - \sum_{i < j} |V_i|$$

Then for $j \in [d] : \mathbf{h}^j = \sigma(A^j \mathbf{h}^{j-1})$


```

import numpy as np

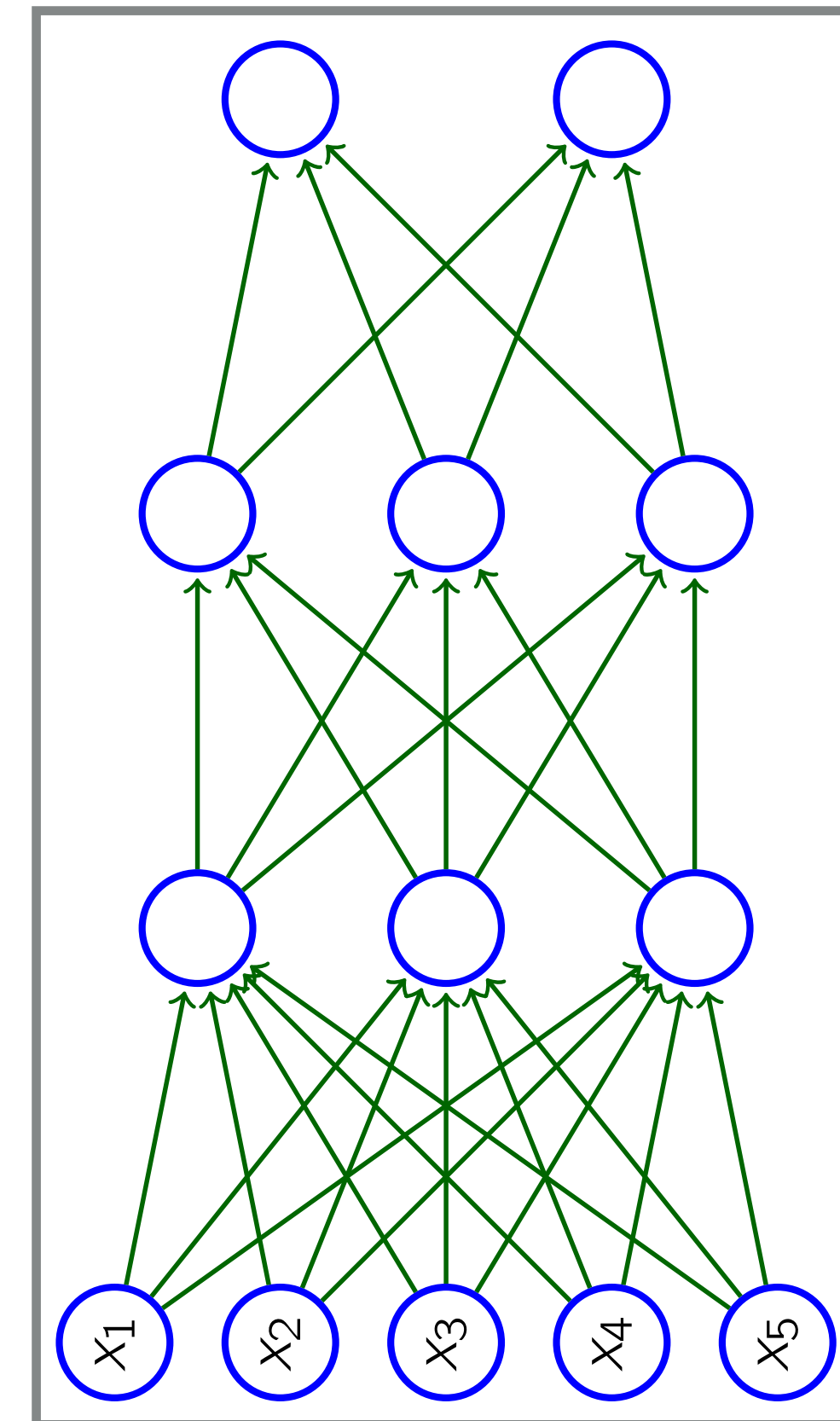
def rand_matrix(m, n):
    return np.random.randn(m, n) / np.sqrt(m)

def MLP(dims):
    As, Bs = [], []
    for i in range(len(dims)-1):
        As.append(rand_matrix(dims[i], dims[i+1]))
        Bs.append(rand_matrix(1, dims[i+1]))
    return (As, Bs)

def predict(net, X, act):
    As, Bs = net
    H = X
    for cd in range(len(As)):
        H = act(H @ As[cd] + Bs[cd])
    return H

dims = [5, 3, 3, 2]
net = MLP(dims)
X = rand_matrix(dims[0], 10).transpose()
O = predict(net, X, np.tanh)

```



Training Neural Networks (Deep Learning)

- MLP defines a non-convex function $f(\mathbf{x})$
- Even if loss $\ell(f(\mathbf{x}), \mathbf{y})$ is convex E.R.M is non-convex
- Initialization is important: symmetry breaking, scale sensitive [more later]
- Gradient-based training takes into account NN's architecture
- Back-propagation: efficient way to calculate gradients using chain rule
$$\nabla_{A_j} \ell(f(\mathbf{x}), \mathbf{y})$$
- Inference & gradient as computations on a graph
- Despite long time to train, often yields good results, many tricks-of-the-trade

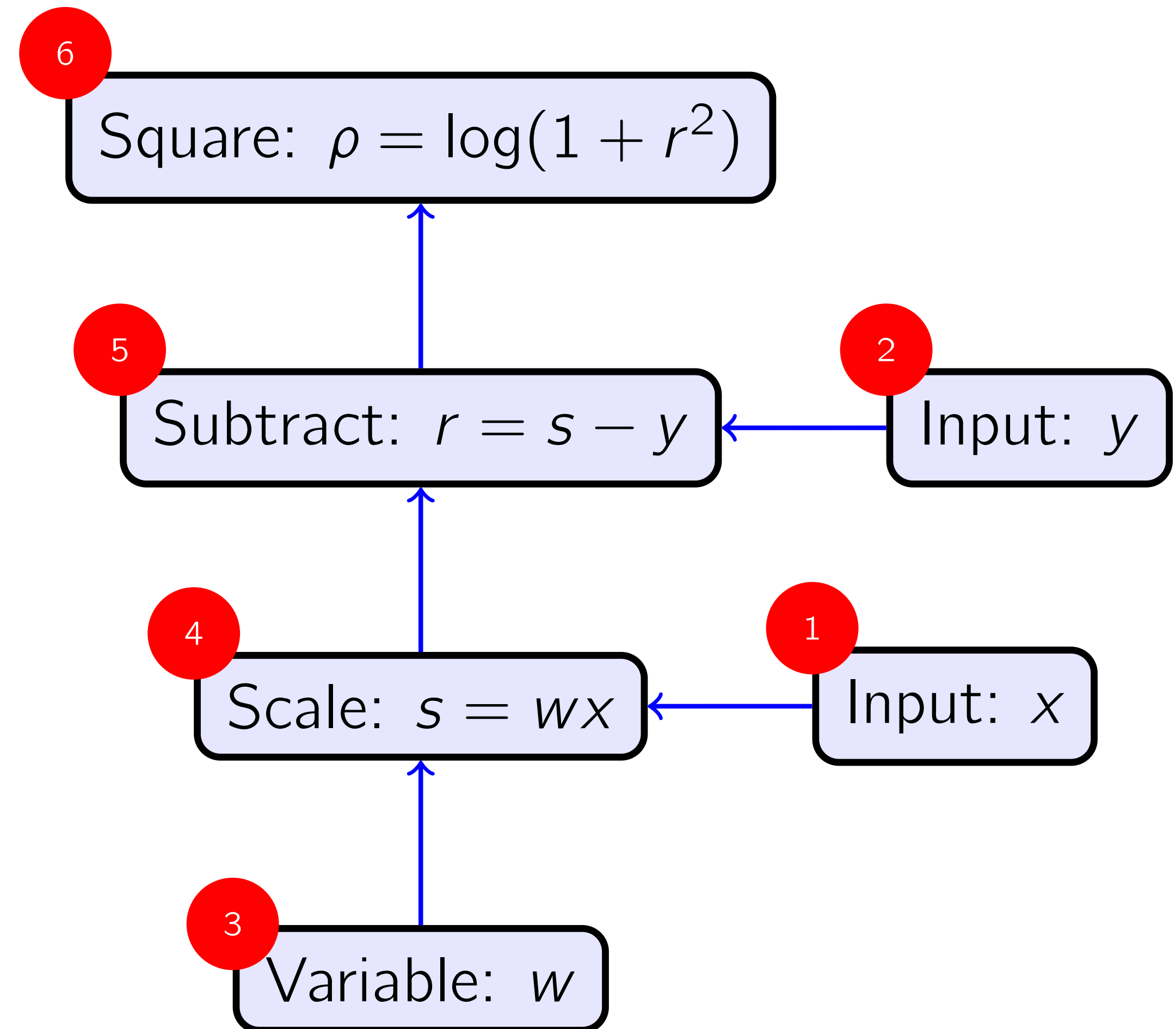
Computation Graph

One dimensional problem

Nodes sorted topologically

Loss: $\ell(y, \hat{y}) = \log(1 + (\hat{y} - y)^2)$

with $\hat{y} = \mathbf{w}\mathbf{x}$ and $\mathbf{x}, \mathbf{y}, \mathbf{w} \in \mathbf{R}$



Derivative Using Chain Rule

$$\rho(z) = \log(1 + z^2)$$

Fix x, y and define functions:

$$r_y(z) = z - y$$

$$s_x(z) = xz$$

Write ℓ as a function of w :

$$\ell(w) = \rho(r_y(s_x(w))) = (\rho \circ r_y \circ s_x)(w)$$

$$\ell'(w) = (\rho \circ r_y \circ s_x)'(w)$$

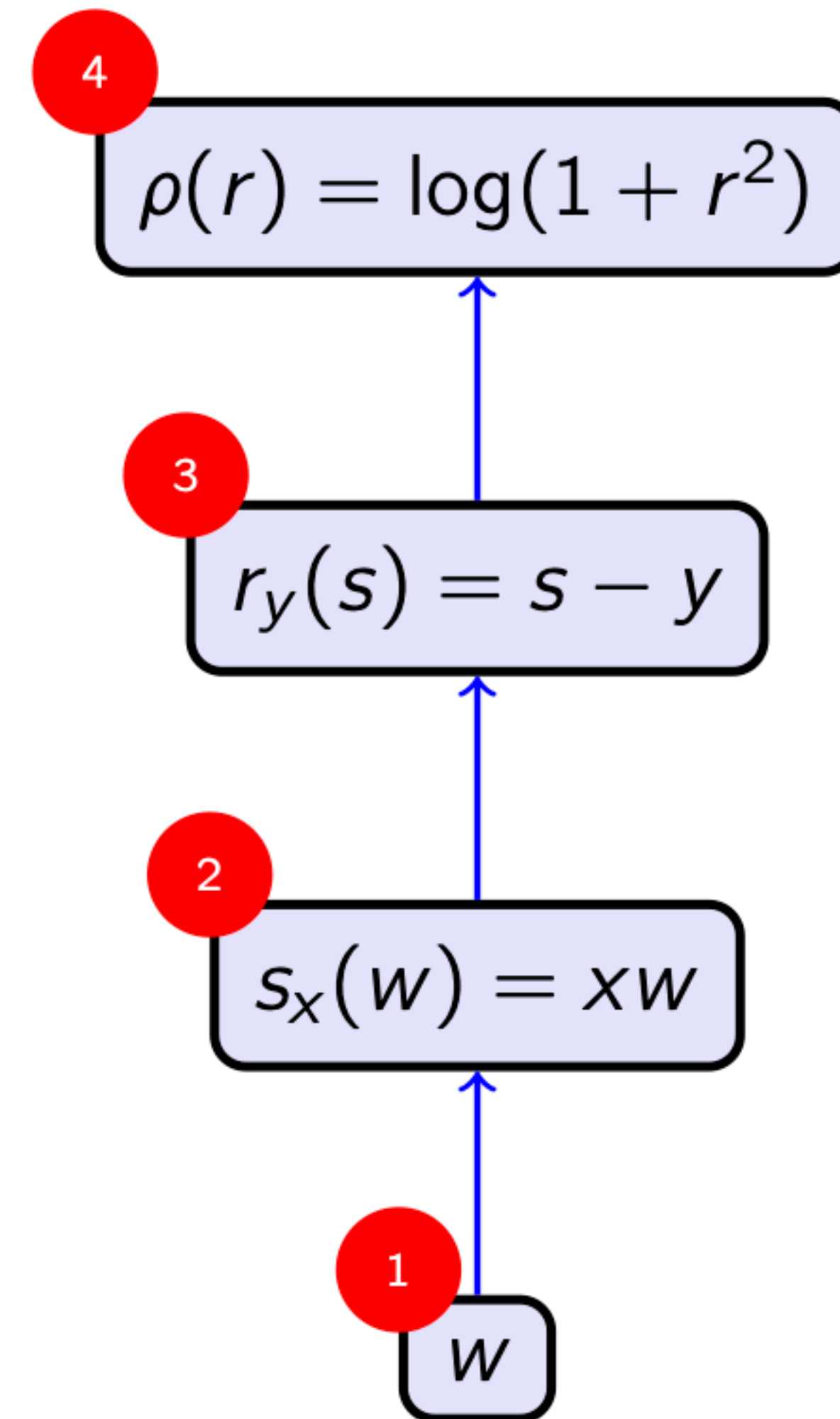
Chain rule:

$$= \rho'(r_y(s_x(w))) \cdot (r_y \circ s_x)'(w)$$

$$= \rho'(r_y(s_x(w))) \cdot r_y'(s_x(w)) \cdot s_x'(w)$$

Inference: Forward Pass

For $v = 1, \dots, m$:
 $v.\text{output} = v.\text{op}(v.\text{inputs}())$

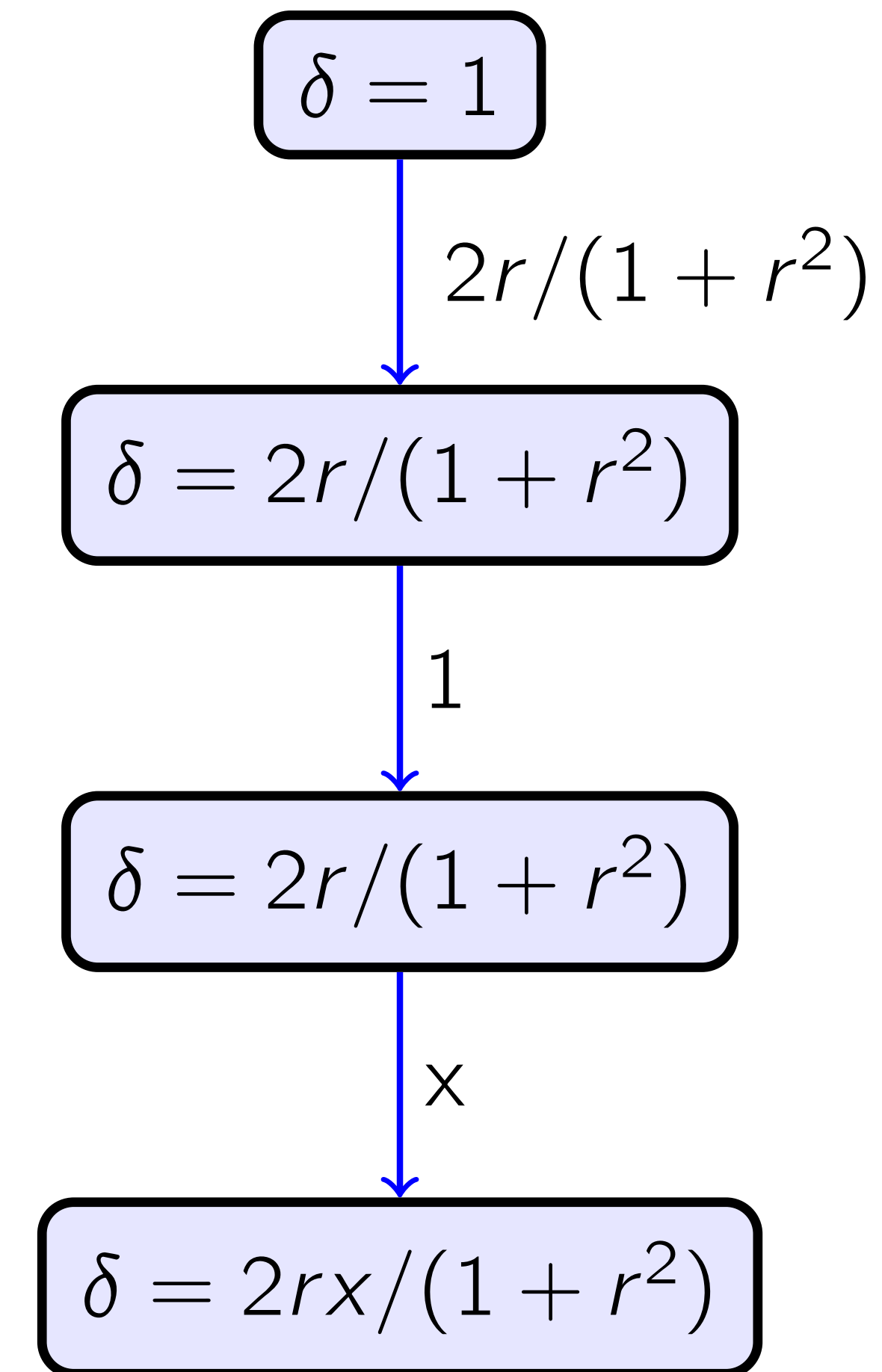


Gradient: Backward Pass (BackProp)

m->delta = 1

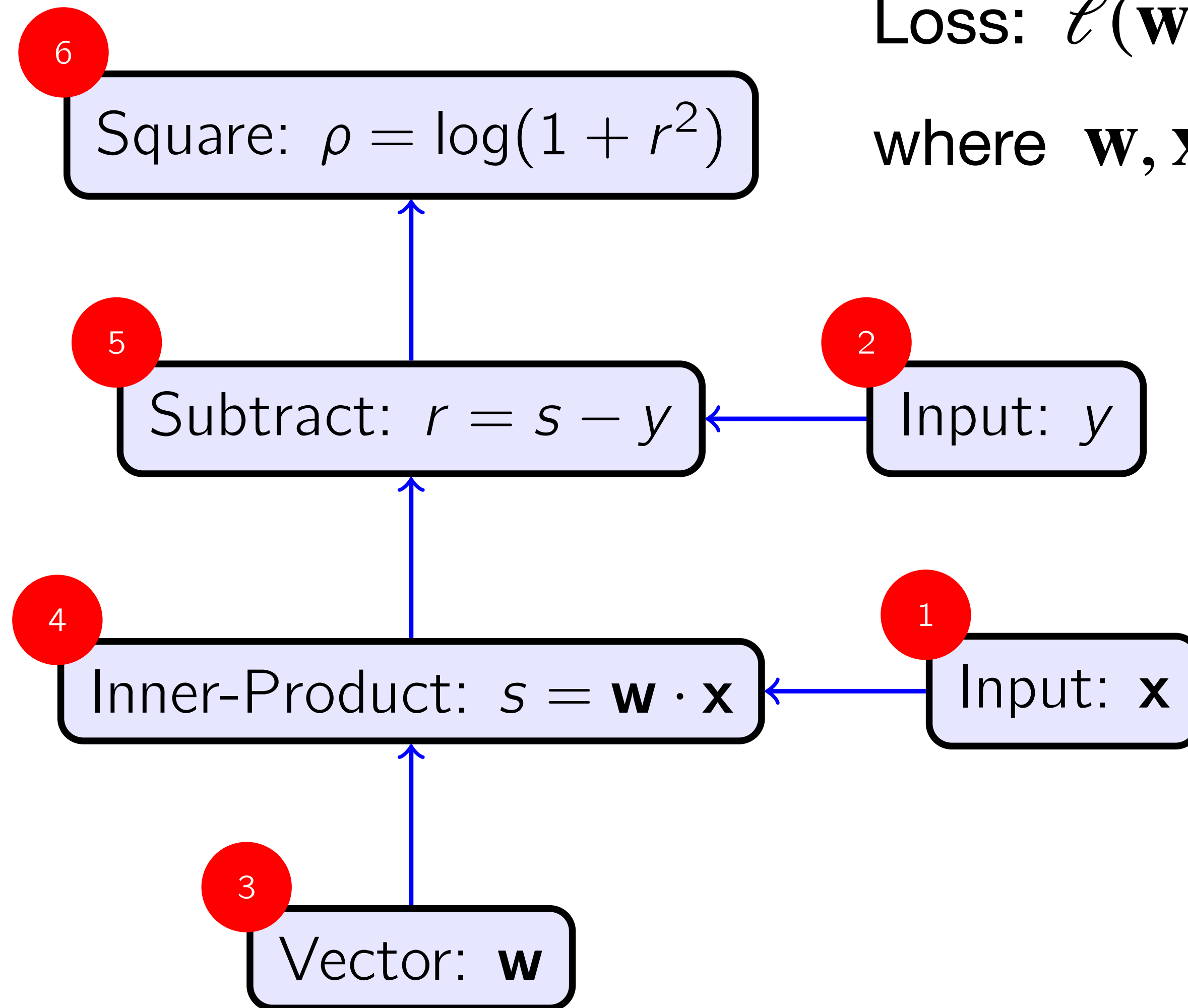
Foreach u s.t. $u \rightarrow v \in E$:

u->delta = v->delta * v->deriv(u)

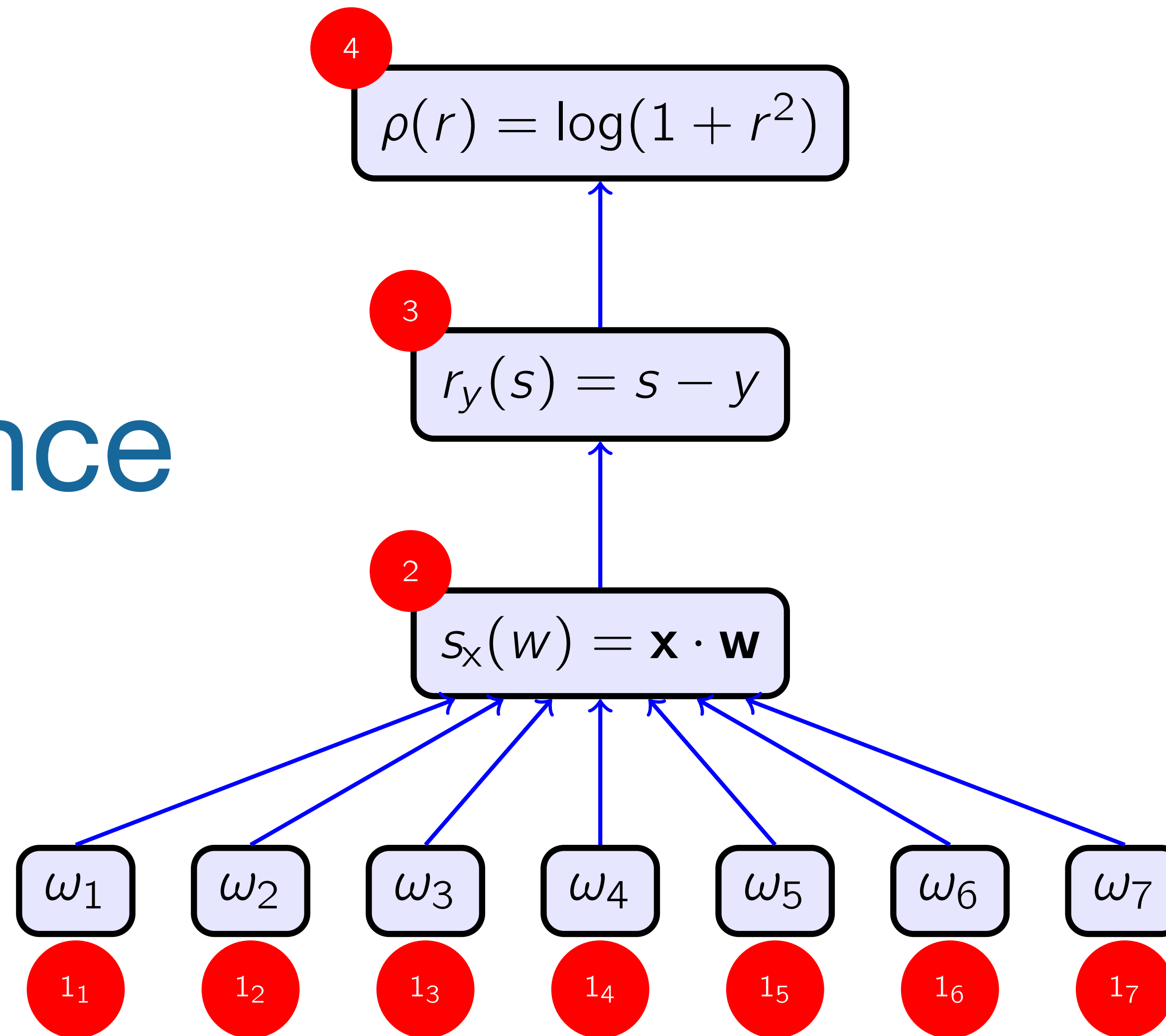


From Scalars to Neurons

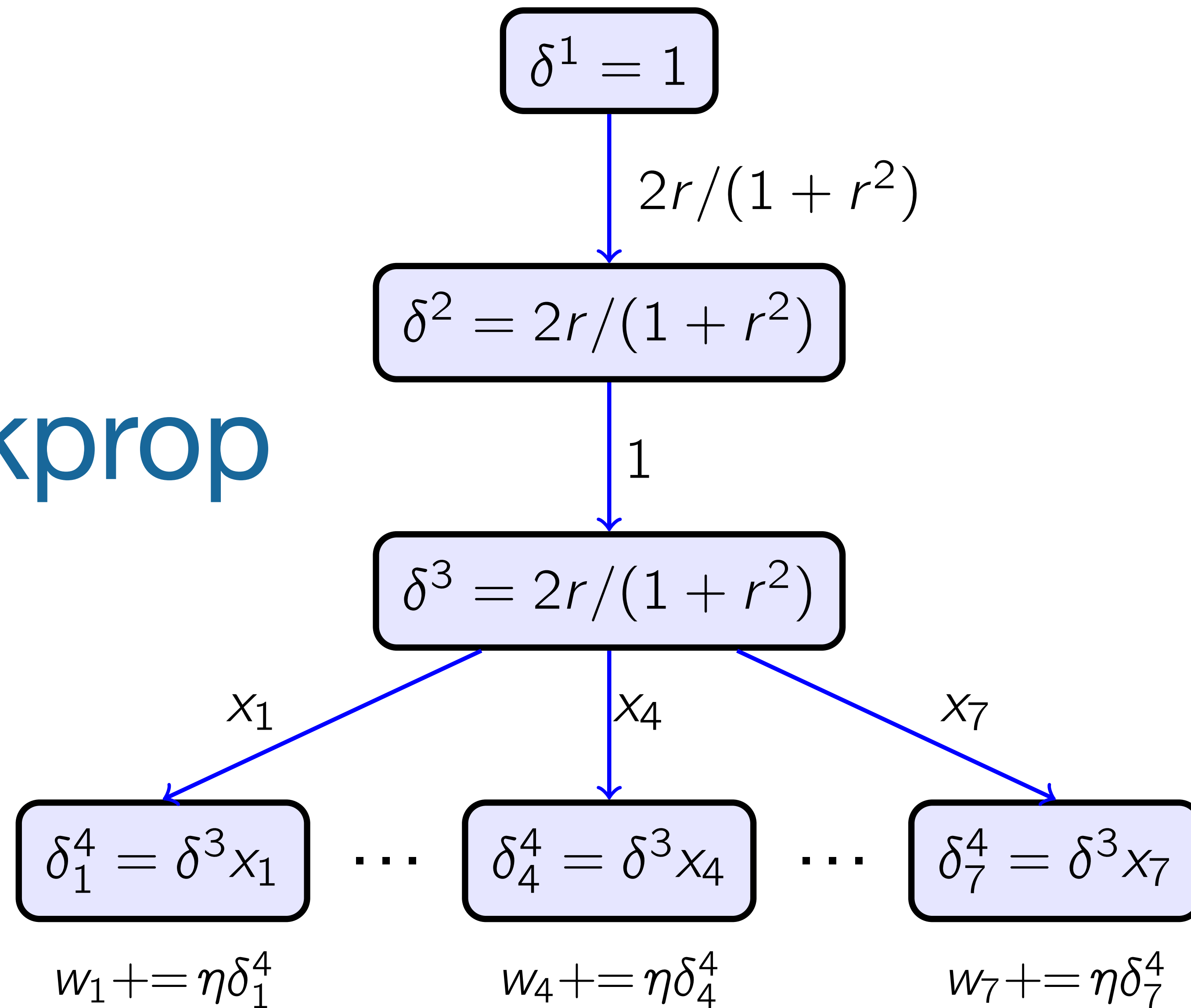
Loss: $\ell(\mathbf{w}, (\mathbf{x}, y)) = \log(1 + (\mathbf{w} \cdot \mathbf{x} - y)^2)$
where $\mathbf{w}, \mathbf{x} \in \mathbf{R}^d$ $y \in \mathbf{R}$



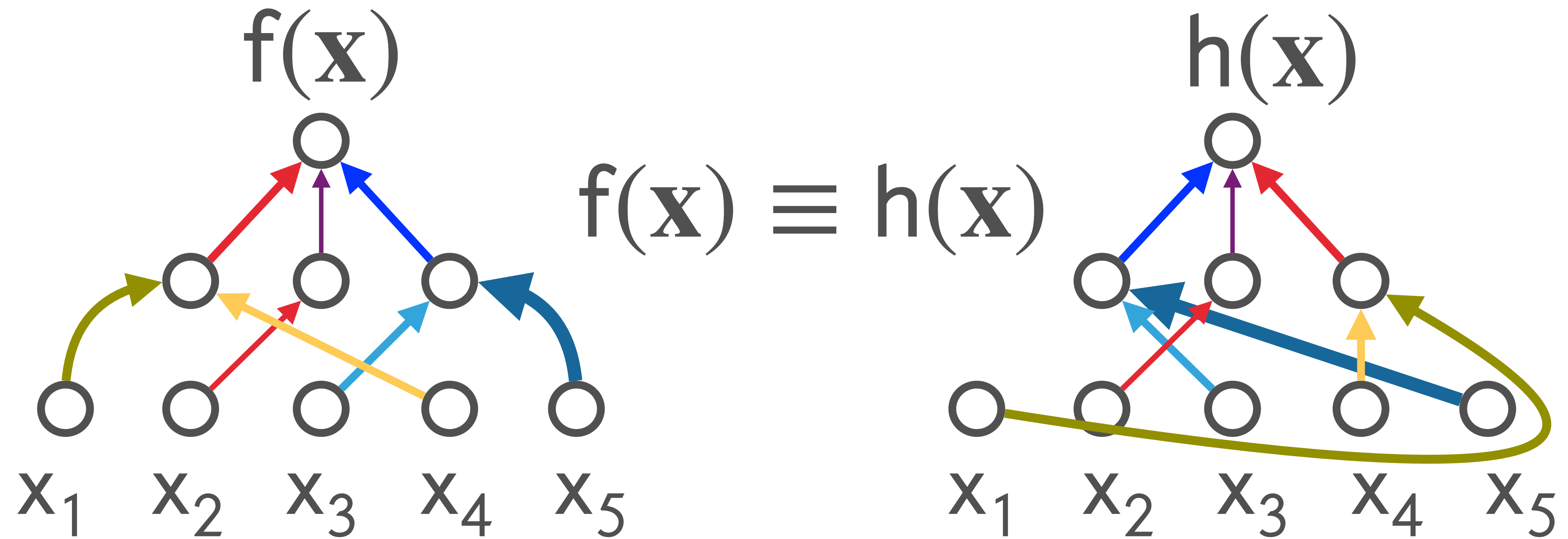
Inference



Backprop



Symmetries



- Input presented in the same order
- Permuting edges and nodes creates seemingly different graphs which are functionally the same

Degrees of Freedom

- ▶ Assume $\sigma(z) = [z]_+$

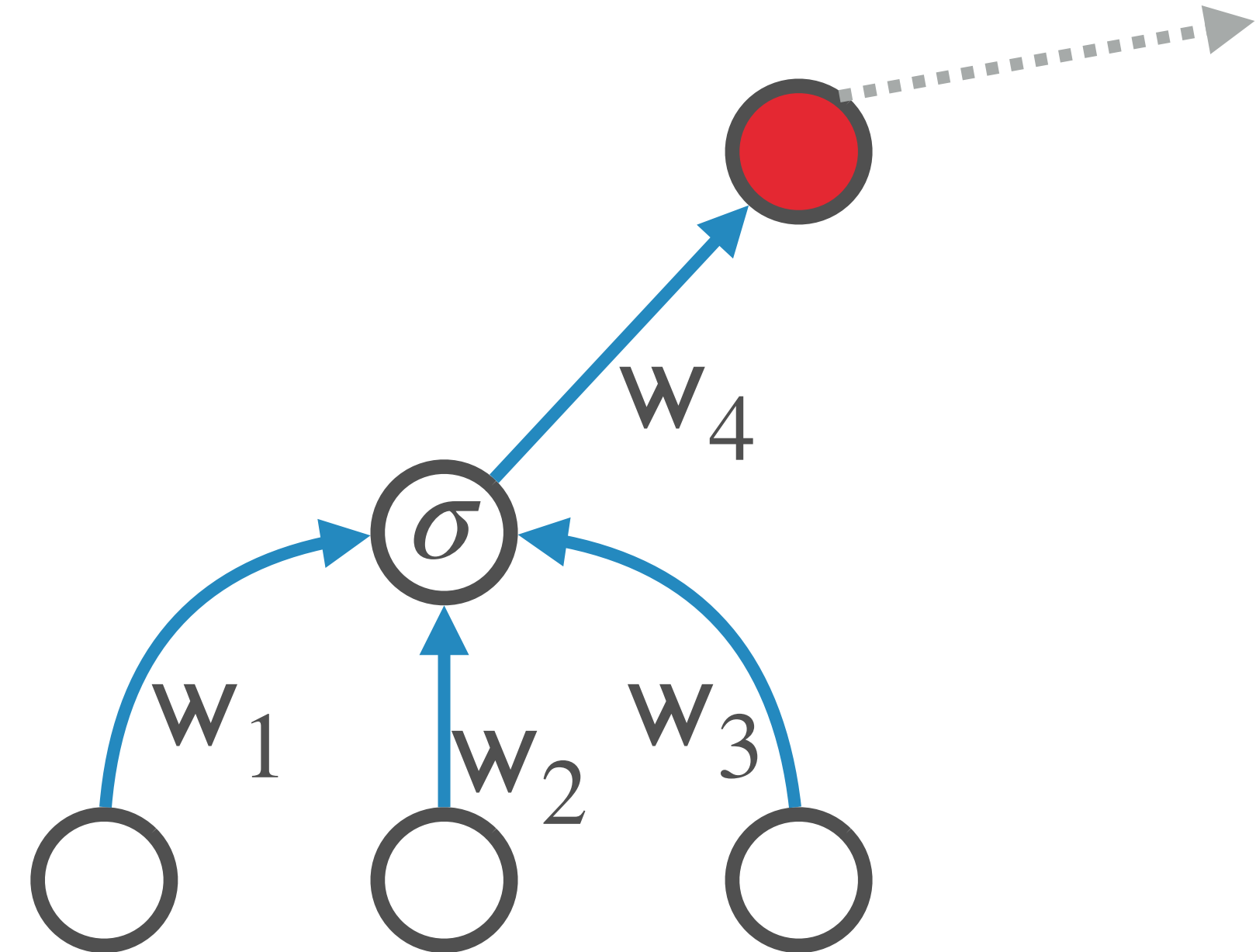
- ▶ Input to **red-neuron**:

$$w_4 \sigma(w_1 x_1 + w_2 x_2 + w_3 x_3)$$

- ▶ Scale : $w_4 \mapsto \frac{w_4}{\alpha}$ and also

$$w_1 \mapsto \alpha w_1 \quad w_2 \mapsto \alpha w_2 \quad w_3 \mapsto \alpha w_3$$

- ▶ Input to **red-neuron** remains the same

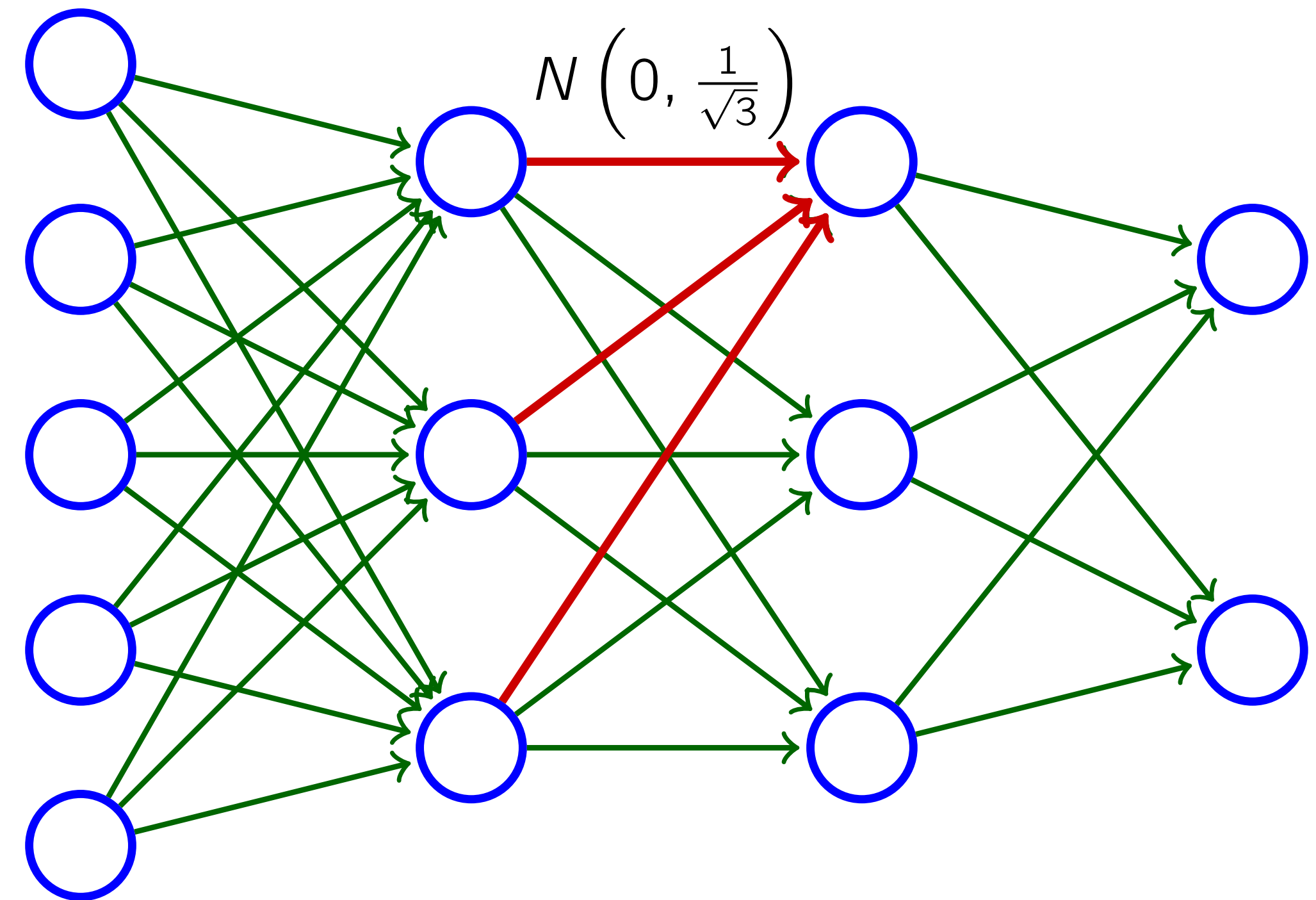


Weight Initialization

To break symmetry and d.o.f.

$$w[u \rightarrow v] \sim \mathcal{N}\left(0, \text{in}(v)^{-\frac{1}{2}}\right)$$

$$\text{where } \text{in}(v) = \left| \{u : u \rightarrow v \in E\} \right|$$



Design an MLP architecture suitable for problem:

- ▶ Number of layers
- ▶ Hidden neurons in each layer
- ▶ Connectivity between layers
- ▶ Activation function

Define loss function between y and $\hat{y} = \text{MLP}(\mathbf{x})$

Initialize weights of MLP

Loop:

Obtain a mini-batch of examples

Perform **Inference** for each example

Perform **Backprop** for each example

Calculate average (over mini-batch) gradient

Update weights