# Backpropagation
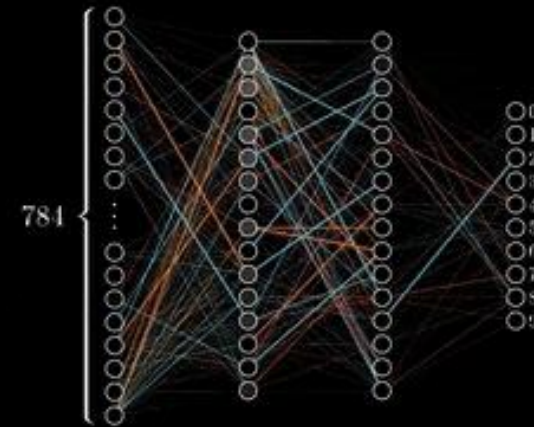
Training in progress. . .

# Backpropagation $\Longleftrightarrow$ Chain rule

○ The neural network loss is a composite function of modules

○ We want the gradient w.r.t. to the parameters of the $l$ layer

$$\frac{d\mathcal{L}}{dw_l} = \frac{d\mathcal{L}}{dh_L} \cdot \frac{dh_L}{dh_{L-1}} \cdot \ldots \cdot \frac{dh_l}{dw_l} \qquad \Rightarrow \qquad \frac{d\mathcal{L}}{dw_l} = \frac{d\mathcal{L}}{dh_l} \cdot \frac{dh_l}{dw_l}$$

Gradient of loss w.r.t. the module output    Gradient of a module w.r.t. its parameters

○ Backpropagation is the "algorithmic manifestation" of the chain rule

# Backpropagation ⟺ Chain rule!!!

o Backpropagating gradients means repeating computation of 2 quantities

$$\frac{d\mathcal{L}}{dw_l} = \frac{d\mathcal{L}}{dh_l} \cdot \frac{dh_l}{dw_l}$$

o For $\frac{dh_l}{dw_l}$ just compute the Jacobian of the $l$-th module w.r.t. to its parameters $w_l$

o <u>Very local</u> rule → "every module looks for its own"

o Since computations can be very local, this means that
  ◦ graphs can be complex
  ◦ modules can be complex if differentiable

# Backpropagation ⟺ Chain rule!!!

○ Backpropagating gradients means repeating computation of 2 quantities

$$\frac{d\mathcal{L}}{dw_l} = \frac{d\mathcal{L}}{dh_l} \cdot \frac{dh_l}{dw_l}$$

○ For $\frac{d\mathcal{L}}{dh_l}$ we apply chain rule again to recursively reuse computations

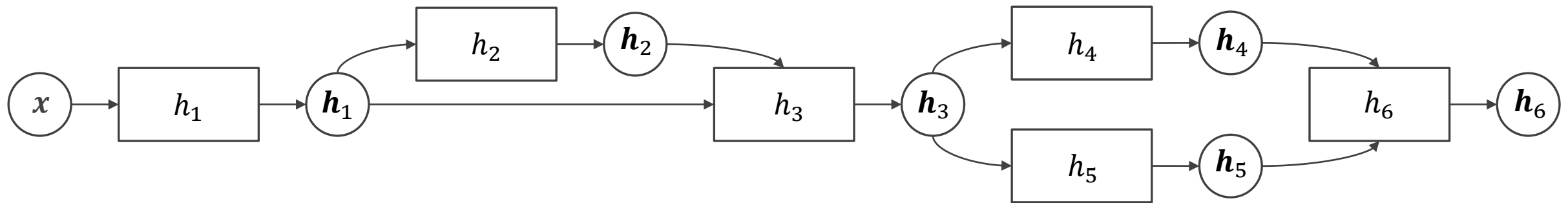$$\frac{d\mathcal{L}}{dh^l} = \frac{d\mathcal{L}}{dh_{l+1}} \cdot \frac{dh_{l+1}}{dh_l}$$

Recursive rule → computation-friendly     Gradient of module w.r.t. its module input

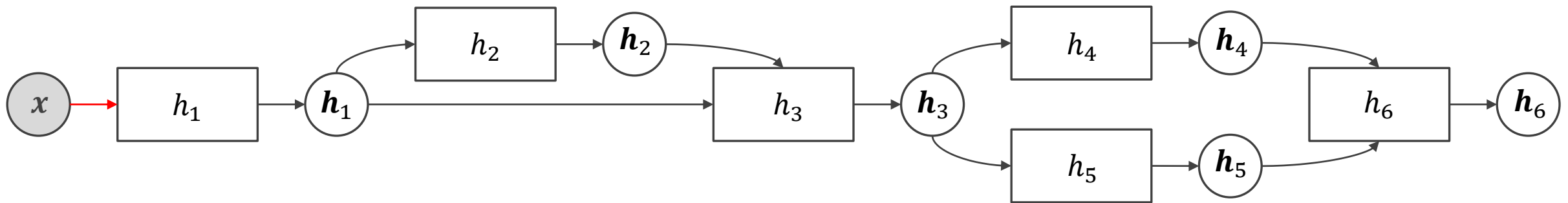○ Remember, the output of a module is the input for the next one: $a_l = x_{l+1}$

# Computational graphs: Forward graph

- Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

- Then, set $x_{l+1} := h_l$

- Store intermediate variables $h_l$
  - will be needed for the backpropagation and saves time at the cost of memory

- Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

o Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

o Then, set $x_{l+1} := h_l$

o Store intermediate variables $h_l$
  ◦ will be needed for the backpropagation and saves time at the cost of memory
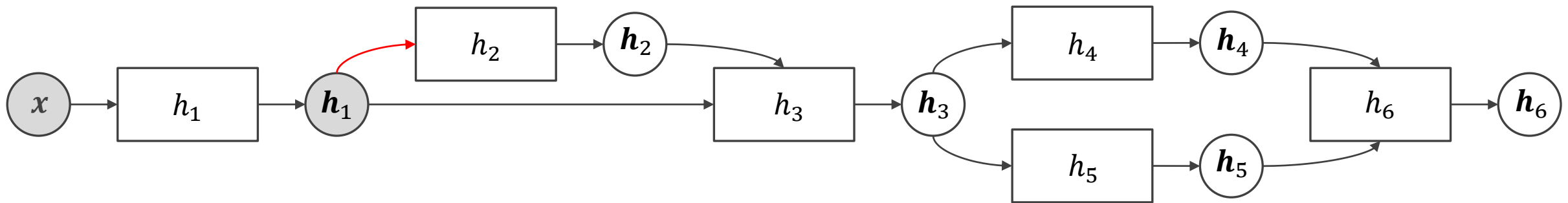
o Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

- Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

- Then, set $x_{l+1} := h_l$

- Store intermediate variables $h_l$
  - will be needed for the backpropagation and saves time at the cost of memory
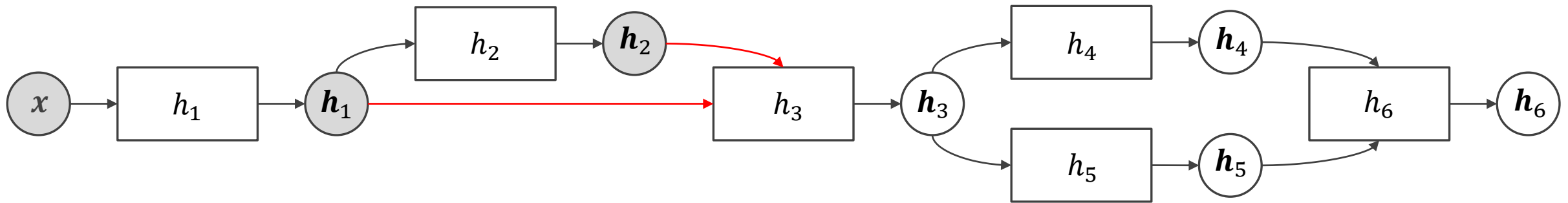
- Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

o Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

o Then, set $x_{l+1} := h_l$

o Store intermediate variables $h_l$
  ◦ will be needed for the backpropagation and saves time at the cost of memory

o Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

o Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

o Then, set $x_{l+1} := h_l$

o Store intermediate variables $h_l$
  ◦ will be needed for the backpropagation and saves time at the cost of memory
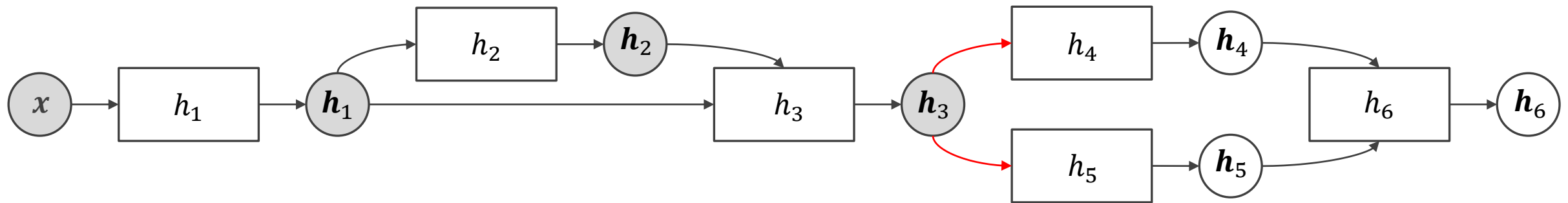
o Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

o Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

o Then, set $x_{l+1} := h_l$

o Store intermediate variables $h_l$
  ◦ will be needed for the backpropagation and saves time at the cost of memory
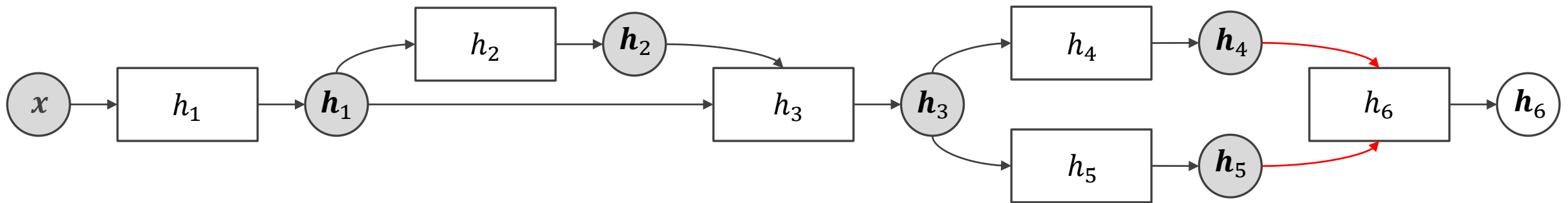
o Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Forward graph

- Compute the activation of each module in the network $\boldsymbol{h}_l = h_l(\boldsymbol{w}; \boldsymbol{x}_l)$

- Then, set $x_{l+1} := h_l$

- Store intermediate variables $h_l$
  - will be needed for the backpropagation and saves time at the cost of memory
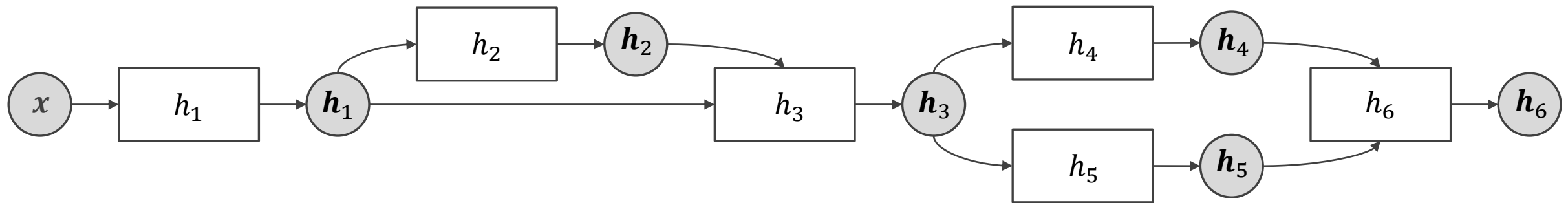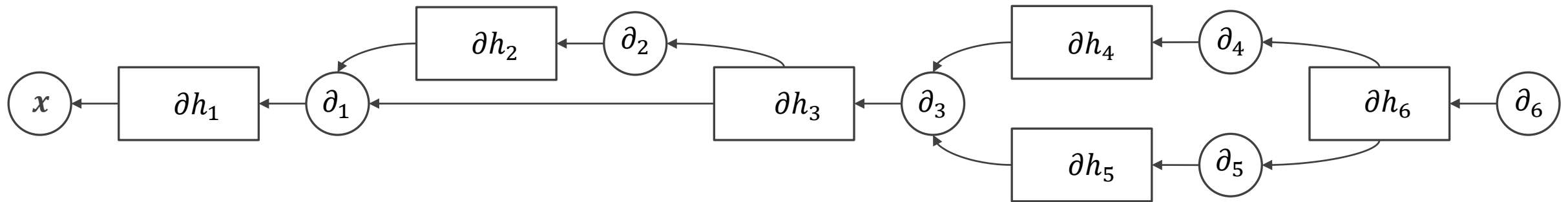
- Then, repeat recursively and in <u>the right order</u>

# Computational graphs: Reverse graph

o Go backwards and use gradient functions instead of activations

◦ Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h^l}{\partial h^{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

o The gradients will need activations from forward propagation, better save them
◦ Sum all gradients from all samples in mini-batch

o Process also known as reverse-mode automatic differentiation
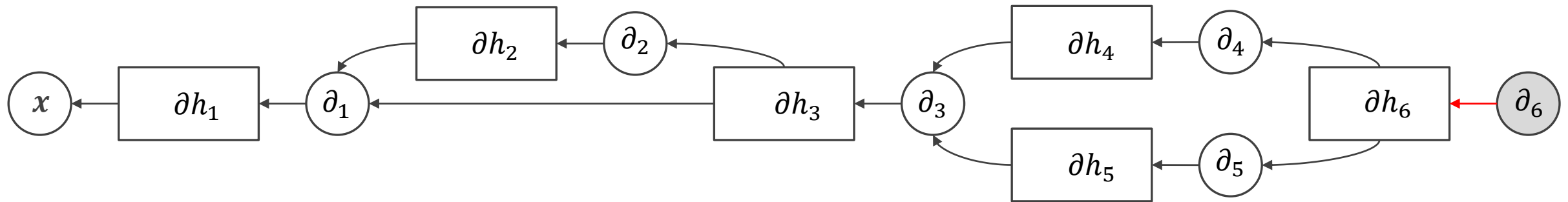◦ Because the flow of computations is reverse to data flow

# Computational graphs: Reverse graph

- Go backwards and use gradient functions instead of activations
  - Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h^l}{\partial h^{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

- The gradients will need activations from forward propagation, better save them
  - Sum all gradients from all samples in mini-batch

- Process also known as reverse-mode automatic differentiation
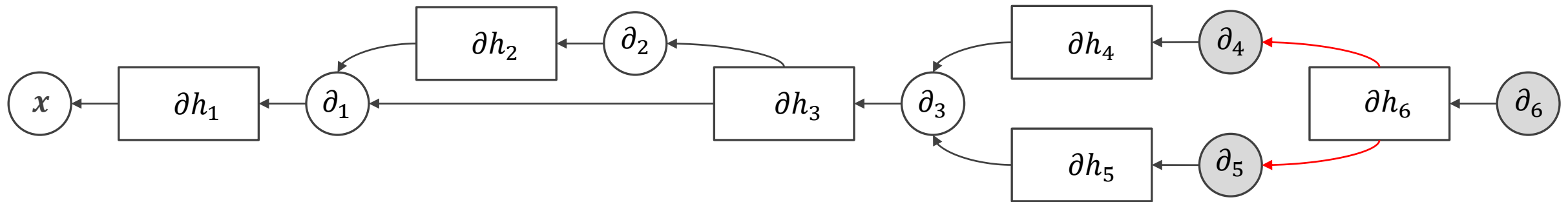  - Because the flow of computations is reverse to data flow

# Computational graphs: Reverse graph

o Go backwards and use gradient functions instead of activations

◦ Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h^l}{\partial h^{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

o The gradients will need activations from forward propagation, better save them
◦ Sum all gradients from all samples in mini-batch

o Process also known as reverse-mode automatic differentiation
◦ Because the flow of computations is reverse to data flow

# Computational graphs: Reverse graph

o Go backwards and use gradient functions instead of activations

   ◦ Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h^l}{\partial h^{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

o The gradients will need activations from forward propagation, better save them
   ◦ Sum all gradients from all samples in mini-batch

o Process also known as reverse-mode automatic differentiation
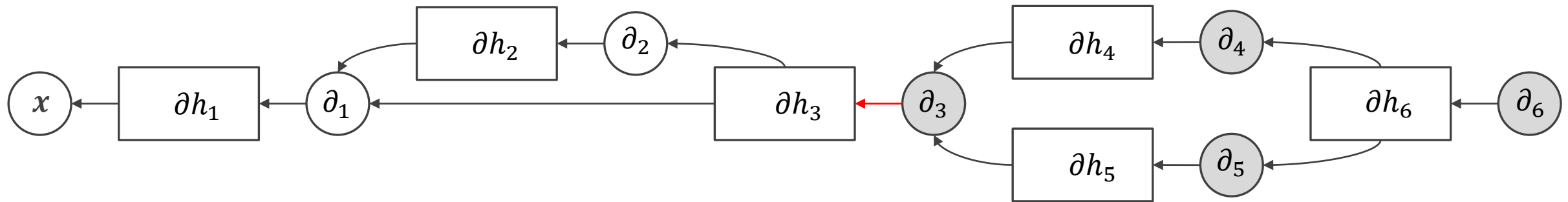   ◦ Because the flow of computations is reverse to data flow

# Computational graphs: Reverse graph

o Go backwards and use gradient functions instead of activations

  ◦ Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h^l}{\partial h^{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

o The gradients will need activations from forward propagation, better save them
  ◦ Sum all gradients from all samples in mini-batch

o Process also known as reverse-mode automatic differentiation
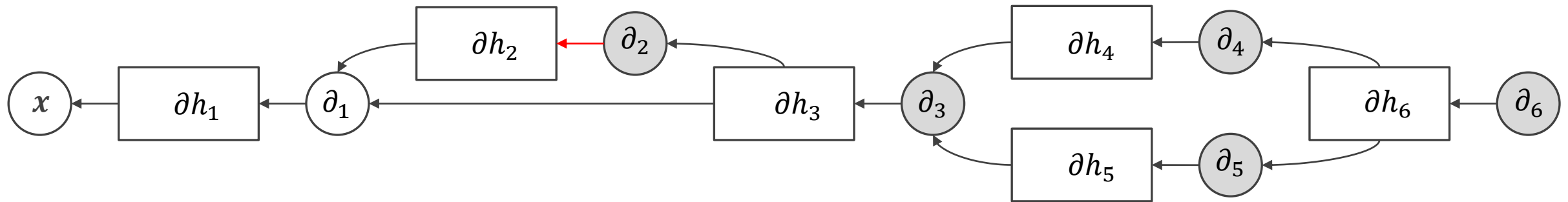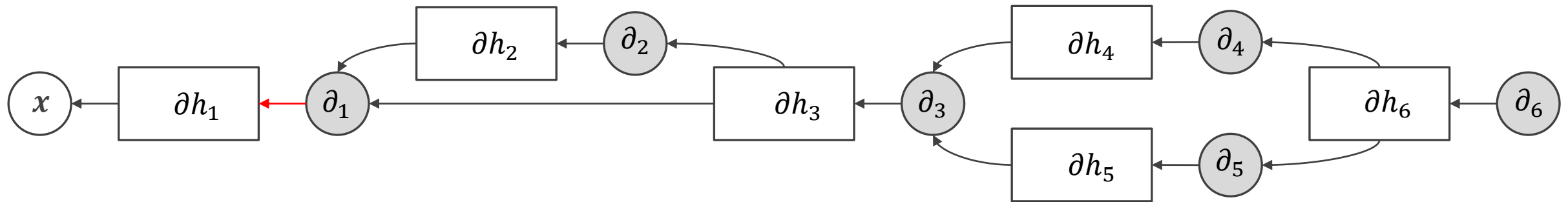  ◦ Because the flow of computations is reverse to data flow

# Computational graphs: Reverse graph

- Go backwards and use gradient functions instead of activations
  - Must have the gradient functions $\frac{\partial h_l}{\partial w_l}, \frac{\partial h_l}{\partial h_{l-1}}$ w.r.t. to $x_l$ & $w_l$ implemented

- The gradients will need activations from forward propagation, better save them
  - Sum all gradients from all samples in mini-batch

- Process also known as reverse-mode automatic differentiation
  - Because the flow of computations is reverse to data flow

# Higher-order derivatives

o Computing higher-order derivatives is similar

$$\text{Hessian} = H = \frac{d^2 h}{d\boldsymbol{x}^2} = \begin{bmatrix} \dfrac{\partial h}{\partial x_1 x_1} & \cdots & \dfrac{\partial h}{\partial x_1 x_M} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial h}{\partial x_M x_M} & \cdots & \dfrac{\partial h}{\partial x_M x_M} \end{bmatrix}$$

o Basically, it is the gradient of the gradient

o Per first-order partial derivative $\frac{\partial h}{\partial x_i}$, auto-differentiate once more

o In practice, computing the second-order gradient is very expensive

# Backprogation in summary

o **Step 1.** Compute forward propagations for all layers recursively

$$h_l = h_l(x_l) \text{ and } x_{l+1} = h_l$$

o **Step 2.** Once done with forward propagation, follow the reverse path.
   ◦ Start from the last layer and for each new layer compute the gradients
   ◦ Cache computations when possible to avoid redundant operations

$$\frac{d\mathcal{L}}{dw_l} = \frac{d\mathcal{L}}{dh_l} \cdot \frac{dh_l}{dw_l} \qquad \frac{d\mathcal{L}}{dh_l} = \frac{d\mathcal{L}}{dh_{l+1}} \cdot \frac{dh_{l+1}}{dh_l}$$

o **Step 3.** Use the gradients $\frac{d\mathcal{L}}{dw^l}$ with Stochastic Gradient Descend to train

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$  → Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$h_2 = \sigma(w_2 h_1)$  → Store $h_2$
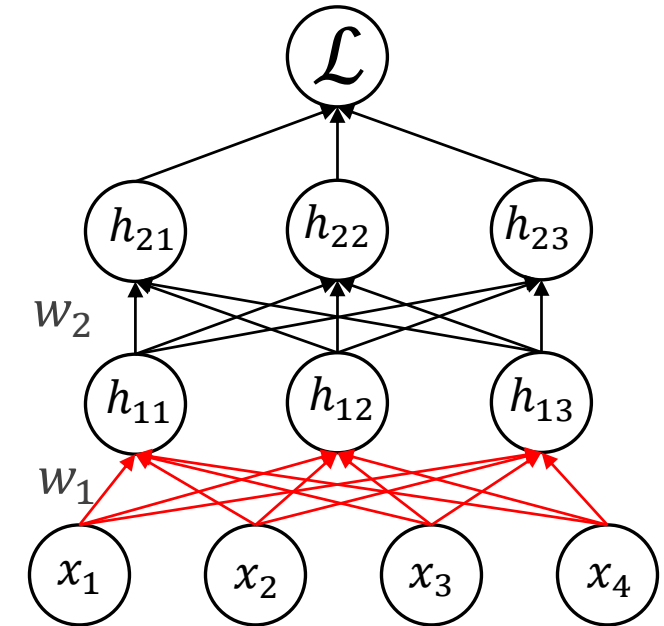
$\mathcal{L} = 0.5 \cdot \| l - h_2 \|^2$

## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1 (1 - h_1)$$

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$ $\qquad \rightarrow$ Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$\rightarrow \quad h_2 = \sigma(w_2 h_1)$ $\qquad \rightarrow$ Store $h_2$
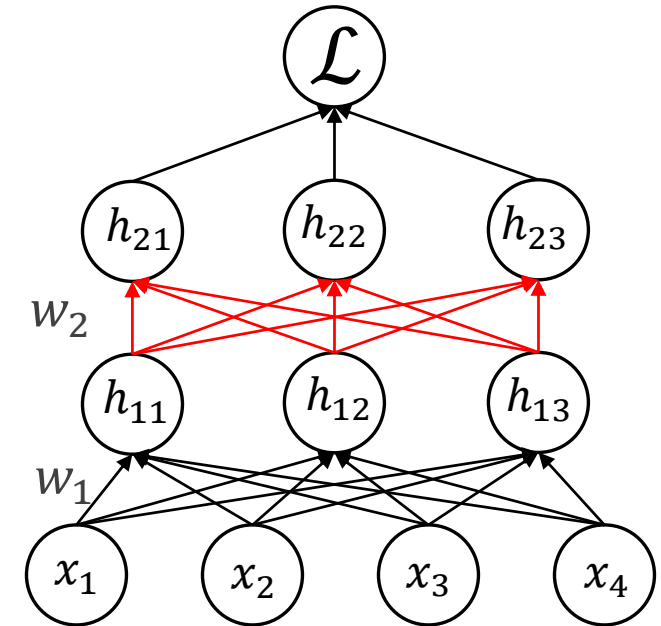
$\mathcal{L} = 0.5 \cdot \| l - h_2 \|^2$

## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2(1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2(1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1(1 - h_1)$$

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$ $\qquad \rightarrow$ Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$h_2 = \sigma(w_2 h_1)$ $\qquad \rightarrow$ Store $h_2$

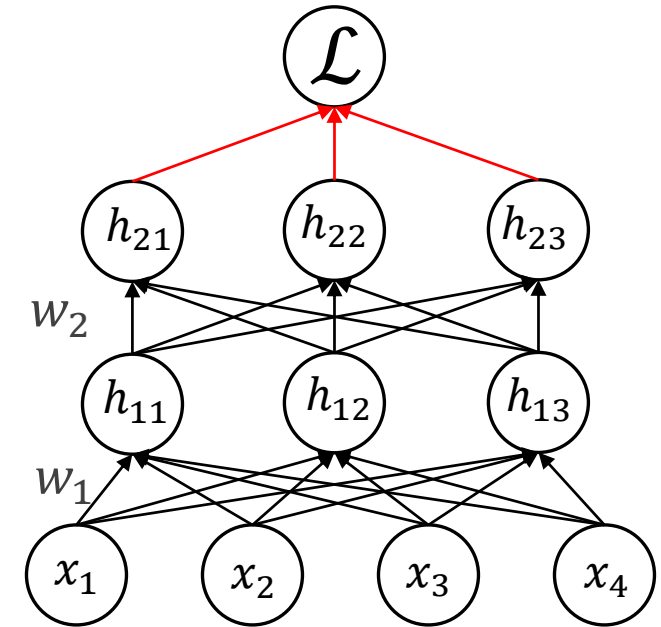$\longrightarrow$ $\mathcal{L} = 0.5 \cdot \| l - h_2 \|^2$

## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1 (1 - h_1)$$

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$        $\rightarrow$ Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$h_2 = \sigma(w_2 h_1)$        $\rightarrow$ Store $h_2$

$\mathcal{L} = 0.5 \cdot \|l - h_2\|^2$

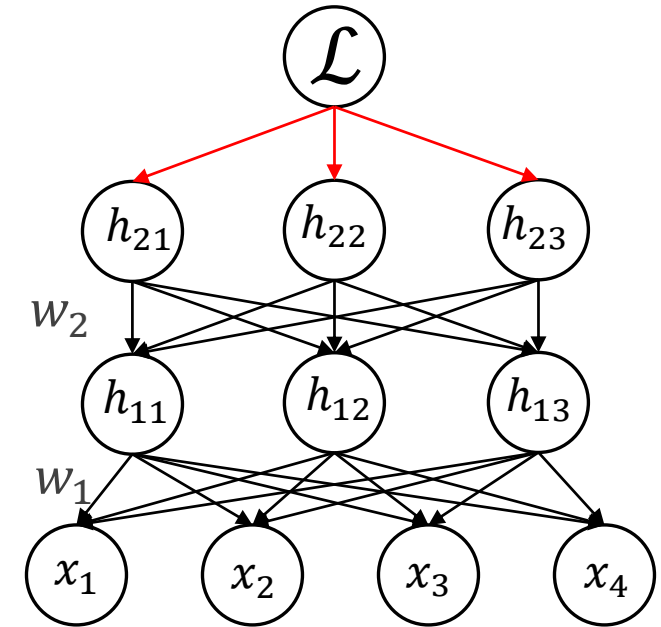## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1 (1 - h_1)$$

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$ $\quad\quad \rightarrow$ Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$h_2 = \sigma(w_2 h_1)$ $\quad\quad \rightarrow$ Store $h_2$
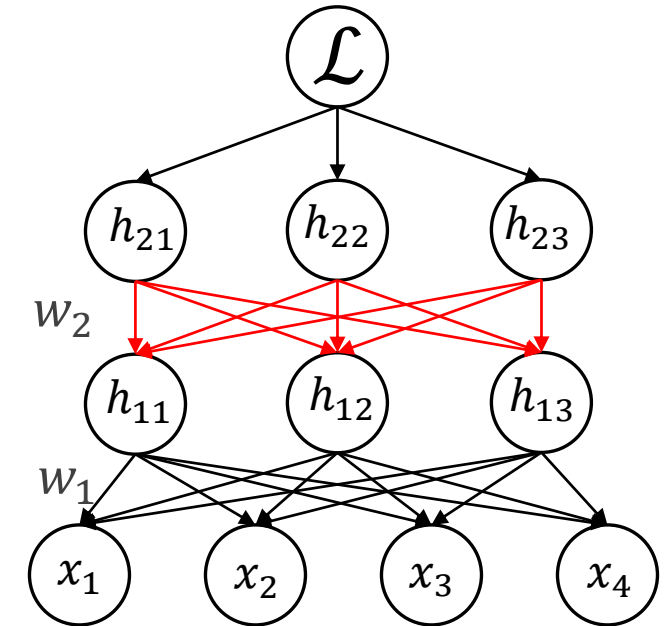
$\mathcal{L} = 0.5 \cdot \|l - h_2\|^2$

## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1 (1 - h_1)$$

# Backpropagation visualization

## Forward propagation

$h_0 = x$

$h_1 = \sigma(w_1 h_0)$      $\rightarrow$ Store $h_1$ . Remember that $\partial_x \sigma = \sigma \cdot (1 - \sigma)$

$h_2 = \sigma(w_2 h_1)$      $\rightarrow$ Store $h_2$
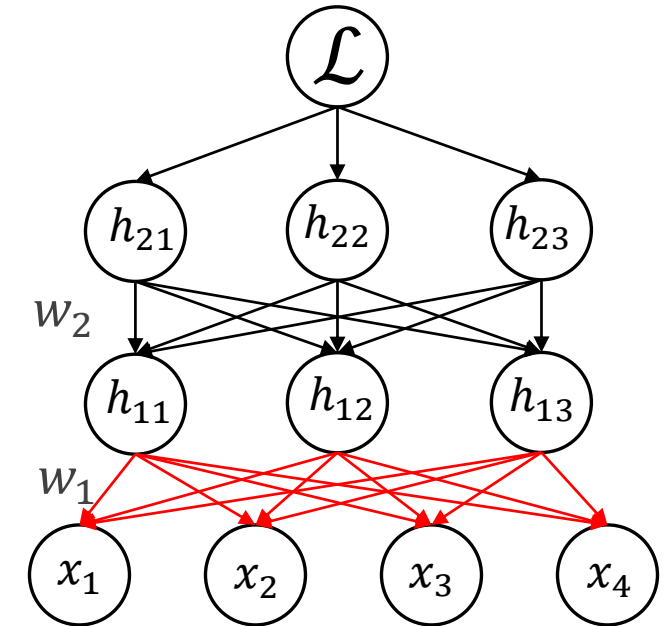
$\mathcal{L} = 0.5 \cdot \|l - h_2\|^2$

## Backward propagation

$$\frac{d\mathcal{L}}{dh_2} = -(y^* - h_2)$$

$$\frac{d\mathcal{L}}{dw_2} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dw_2} = \frac{d\mathcal{L}}{dh_2} h_1 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{da_2} h_1 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dh_1} = \frac{d\mathcal{L}}{dh_2}\frac{dh_2}{dh_1} = \frac{d\mathcal{L}}{dh_2} w_2 \sigma(w_2 h_1)\big(1 - \sigma(w_2 h_1)\big) = \frac{d\mathcal{L}}{dh_2} w_2 h_2 (1 - h_2)$$

$$\frac{d\mathcal{L}}{dw_1} = \frac{d\mathcal{L}}{dh_1}\frac{dh_1}{dw_1} = \frac{d\mathcal{L}}{dh_1} h_0 \sigma(w_1 h_0)\big(1 - \sigma(w_1 h_0)\big) = \frac{d\mathcal{L}}{dh_1} h_0 h_1 (1 - h_1)$$

# What's the big deal?

o Backpropagation is as simple as it is complicated

o Mathematically, just the chain rule

o That simple, that we can even automate it ("reverse-mode differentiation")

o However, algorithmically the devil is in the details to make it efficient

o And, theoretically, why does it even work given the strong non-convexity?

# Summary

o Modularity in Neural Networks

o Neural Network Modules

o Neural Network Cheatsheet

o Backpropagation

**Reading material**

o Chapter 6

o Efficient Backprop, LeCun et al., 1998