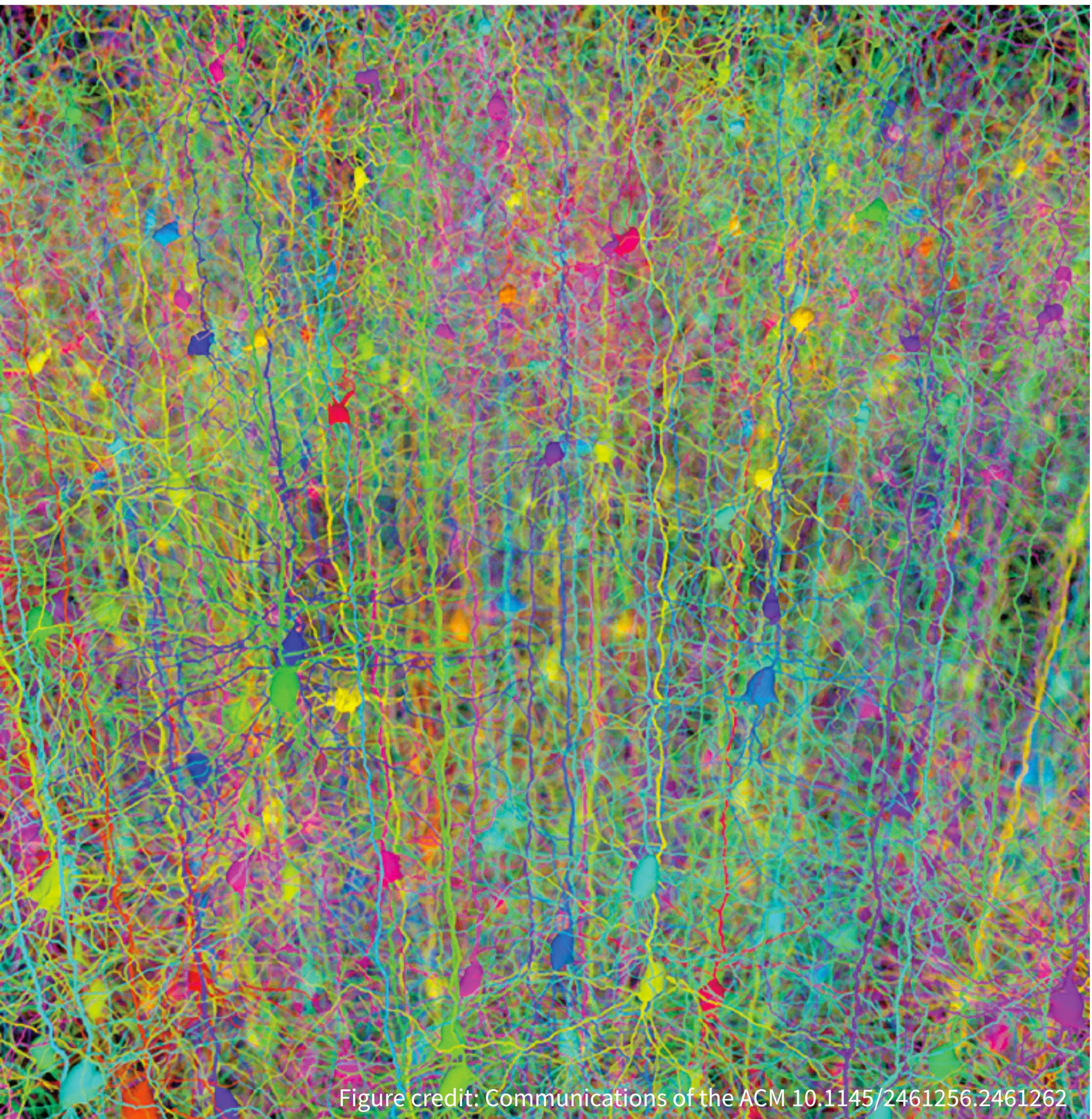


# Training FeedForward Networks

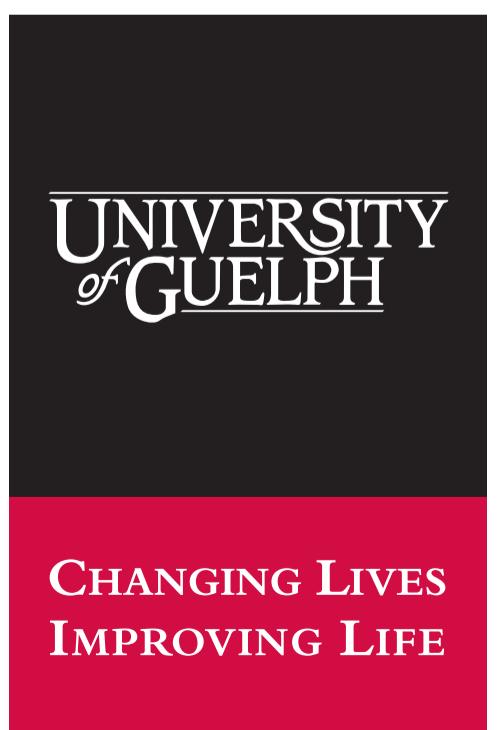


GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING  
UNIVERSITY OF GUELPH

CANADIAN INSTITUTE  
FOR ADVANCED RESEARCH



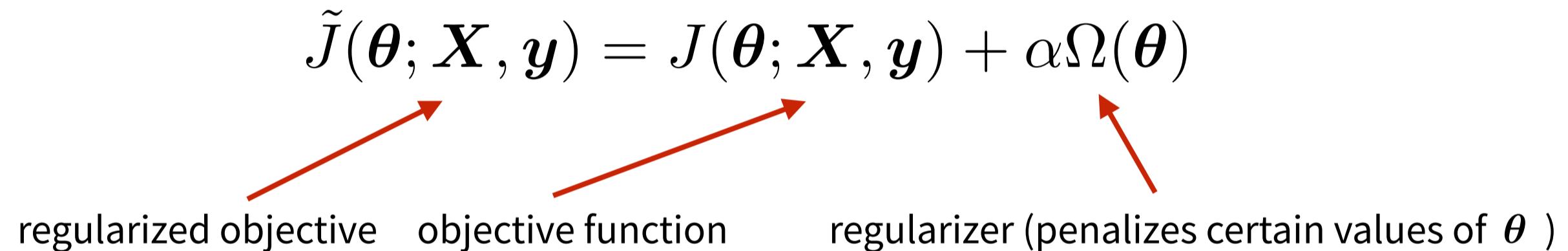
**CIFAR**  
CANADIAN  
INSTITUTE  
FOR  
ADVANCED  
RESEARCH

# Empirical Risk Minimization

A framework to design learning algorithms:

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(\theta)$$

regularized objective    objective function    regularizer (penalizes certain values of  $\theta$ )



Learning cast as optimization:

$$\arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m L \left( f \left( \mathbf{x}^{(i)}; \theta \right), y^{(i)} \right) + \alpha\Omega(\theta)$$

Note: loss function is typically  
a surrogate for what we truly should optimize

# Stochastic Gradient Descent

# Stochastic Gradient Descent

An algorithm that performs parameter updates after each example:

# Stochastic Gradient Descent

An algorithm that performs parameter updates after each example:

- initialize:  $\theta \quad (\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(l)}, \mathbf{b}^{(l)}\})$
- while stopping criterion not met, do:
  - for each minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with targets  $y^{(i)}$ 
    - Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) + \alpha \nabla_{\theta} \Omega(\theta)$
    - Apply update:  $\theta \leftarrow \theta - \epsilon g$

# Stochastic Gradient Descent

An algorithm that performs parameter updates after each example:

- initialize:  $\theta \quad (\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(l)}, \mathbf{b}^{(l)}\})$
- while stopping criterion not met, do:
  - for each minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with targets  $y^{(i)}$ 
    - Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) + \alpha \nabla_{\theta} \Omega(\theta)$
    - Apply update:  $\theta \leftarrow \theta - \epsilon g$

To apply this algorithm to neural network training, we need:

# Stochastic Gradient Descent

An algorithm that performs parameter updates after each example:

- initialize:  $\theta \quad (\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(l)}, \mathbf{b}^{(l)}\})$
- while stopping criterion not met, do:
  - for each minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with targets  $y^{(i)}$ 
    - Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) + \alpha \nabla_{\theta} \Omega(\theta)$
    - Apply update:  $\theta \leftarrow \theta - \epsilon g$

To apply this algorithm to neural network training, we need:

- the loss function  $L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
- a procedure to compute parameter gradients  $\nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
- the regularizer  $\Omega(\theta)$  (and the gradient  $\nabla_{\theta} \Omega(\theta)$ )
- an initialization method

# Loss Function

Regression (Mean-squared error):

$$L_{MSE} = \frac{1}{2} (f(\mathbf{x}) - y)^2$$

Classification:

- Let neural network estimate  $f(\mathbf{x})_c = p(y = c|\mathbf{x})$ 
  - We could maximize the probabilities of  $y^{(i)}$  given  $\mathbf{x}^{(i)}$  in the training set
- To frame as minimization, minimize the negative log-likelihood

$$L_{CE} = - \sum_c 1_{(y=c)} \log f(\mathbf{x})_c = - \log f(\mathbf{x})_y$$

“cross-entropy” 

- take the log to simplify numerical stability and mathematical simplicity

# Forward and Back-propagation

- The process of propagating  $x$  through a feedforward neural network to produce an output  $\hat{y}$  is called **forward propagation**
- During training, forward propagation continues to produce a scalar cost  $J(\theta)$
- The **back-propagation** algorithm (backprop) allows the cost to flow backward through the net to compute the gradient

# Misunderstanding?

# Misunderstanding?

- Backprop is **not the whole learning algorithm** for neural networks
- Back-propagation only refers to the method for computing the **gradient**
- Another algorithm, such as stochastic gradient descent (SGD) is used to perform learning, **using this gradient**

# Misunderstanding?

- Backprop is **not the whole learning algorithm** for neural networks
- Back-propagation only refers to the method for computing the **gradient**
- Another algorithm, such as stochastic gradient descent (SGD) is used to perform learning, **using this gradient**

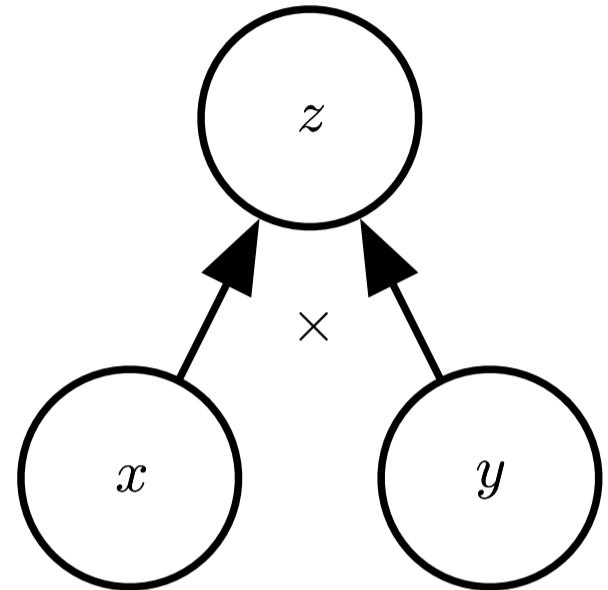
Terminology: backprop computes gradient  $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$

- for an arbitrary function  $f$
- where  $\mathbf{x}$  is a set of variables whose derivatives are desired
- and  $\mathbf{y}$  are inputs to the function whose derivatives are not desired

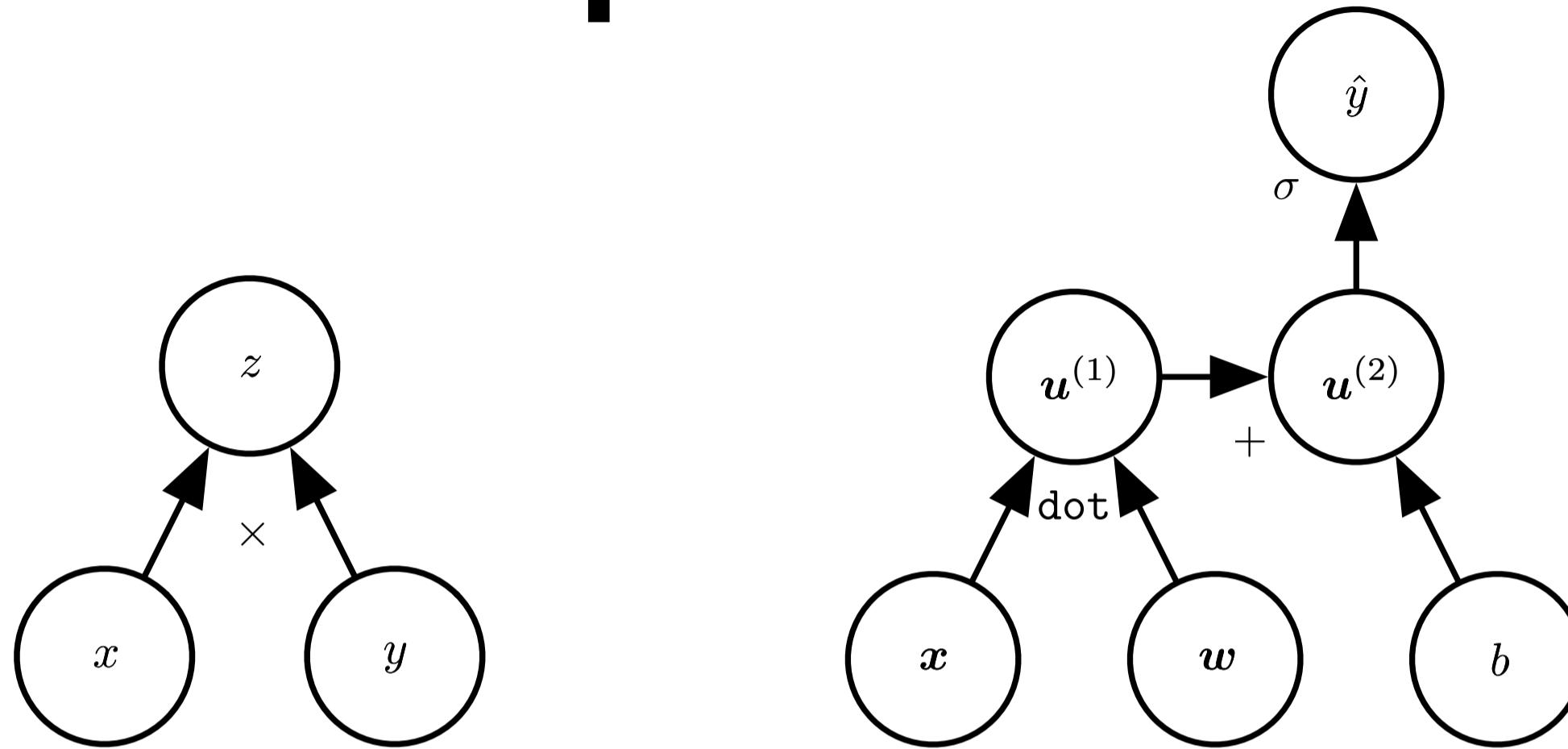
# Operations

- To introduce the **computational graph** language, we need to introduce the idea of an operation
- An **operation** is a simple function of one or more variables
- A graph language is accompanied by a set of allowable operations
- More complicated functions are made by composing many operations
- Without loss of generality, we define an operation to **return only a single variable** (the output can have multiple entries)

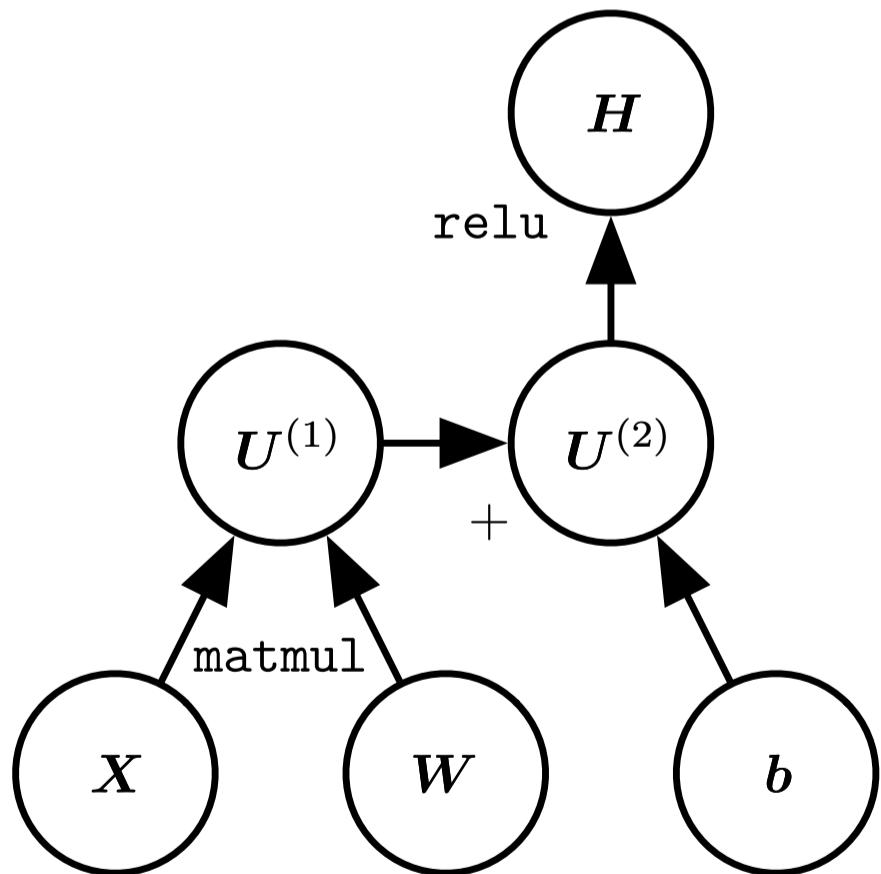
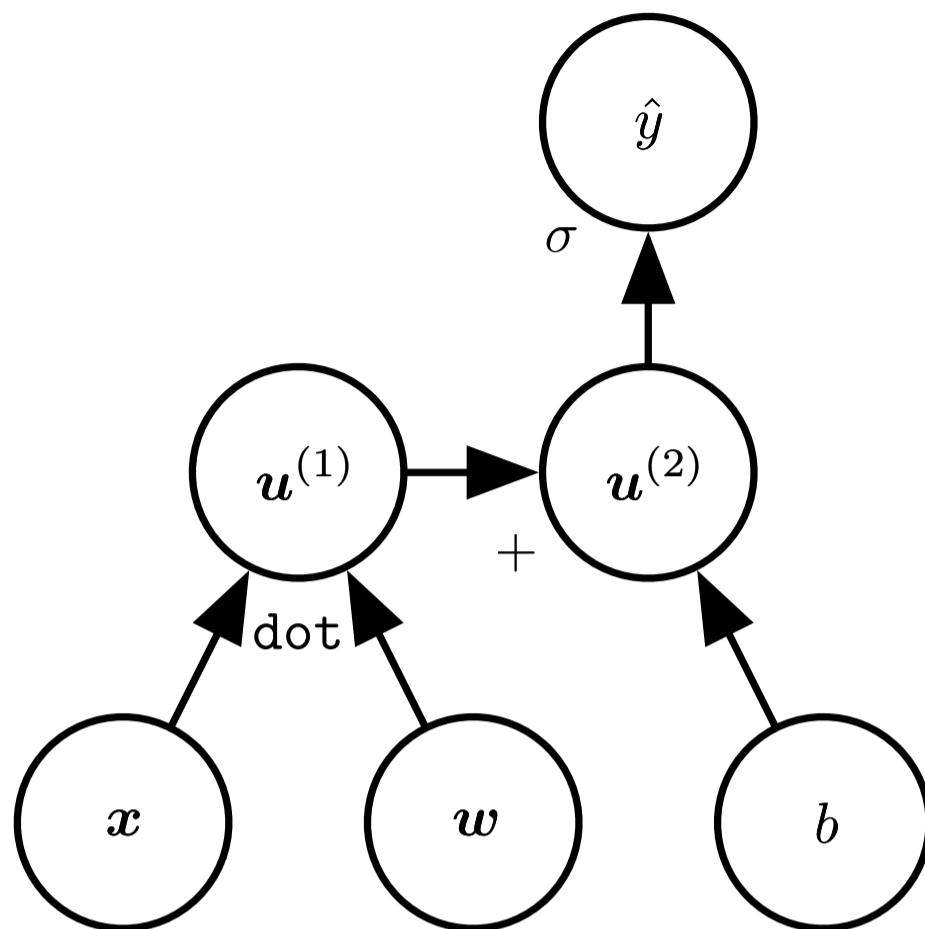
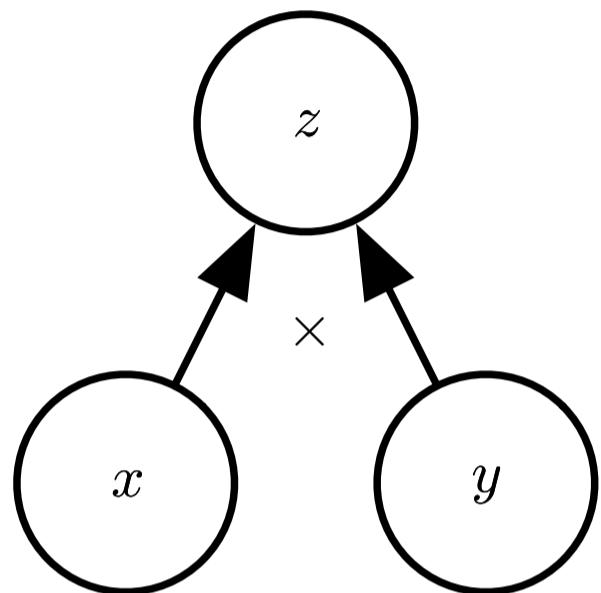
# Example: Computational Graphs



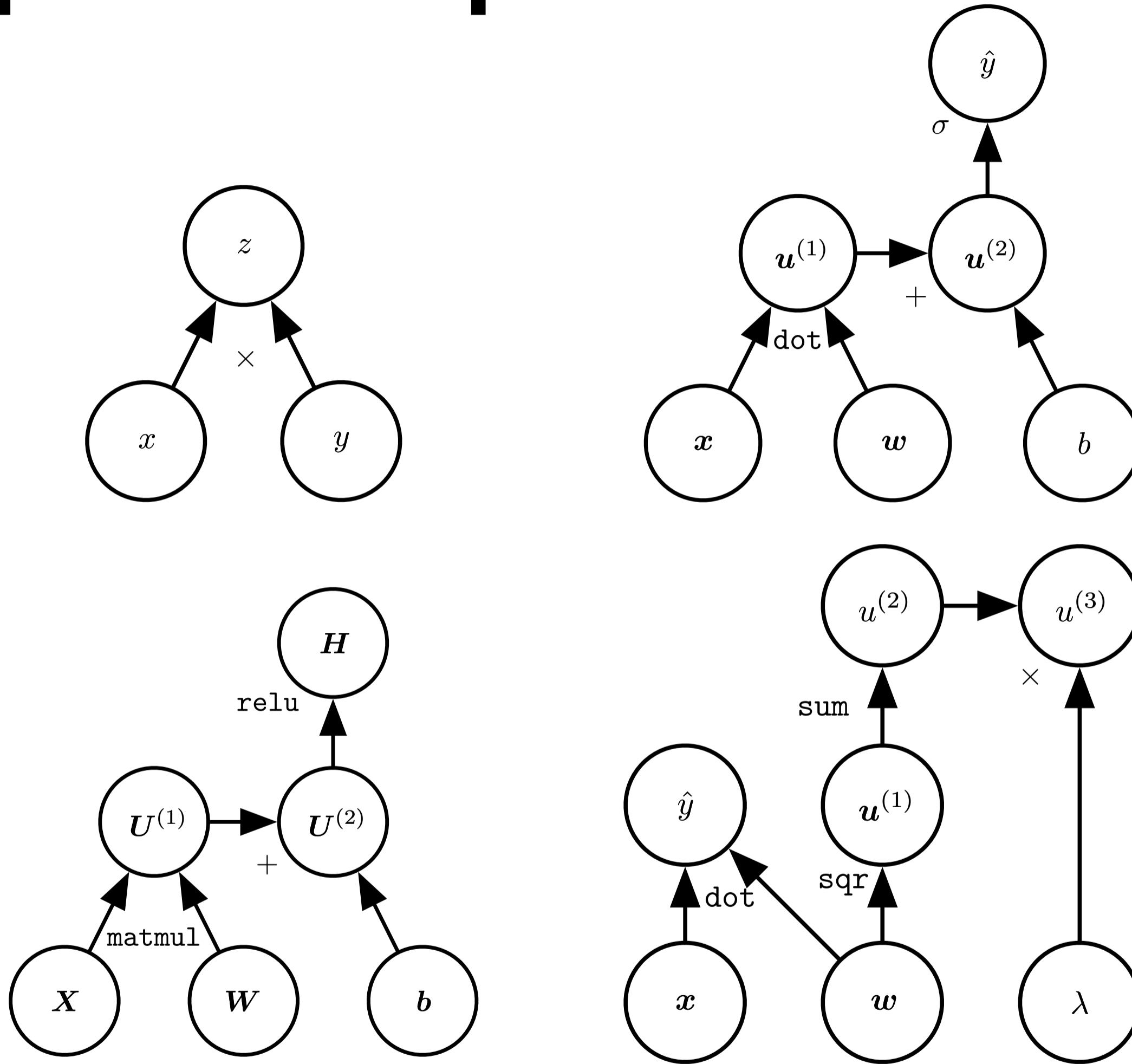
# Example: Computational Graphs



# Example: Computational Graphs



# Example: Computational Graphs



# Chain Rule of Calculus

# Chain Rule of Calculus

Used to compute the derivatives of functions formed by composing other functions whose derivatives are known

- Suppose that  $y = g(x)$  and  $z = f(g(x)) = f(y)$
- Then the chain rule states that:

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

# Chain Rule of Calculus

Used to compute the derivatives of functions formed by composing other functions whose derivatives are known

- Suppose that  $y = g(x)$  and  $z = f(g(x)) = f(y)$
- Then the chain rule states that:

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

We can generalize this beyond the scalar case

- Suppose that  $\mathbf{y} = g(\mathbf{x})$  and  $z = f(\mathbf{y})$ , then  $\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$
- In vector notation, this can be equivalently rewritten as:

$$\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

# Chain Rule of Calculus

Used to compute the derivatives of functions formed by composing other functions whose derivatives are known

- Suppose that  $y = g(x)$  and  $z = f(g(x)) = f(y)$
- Then the chain rule states that:

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

We can generalize this beyond the scalar case

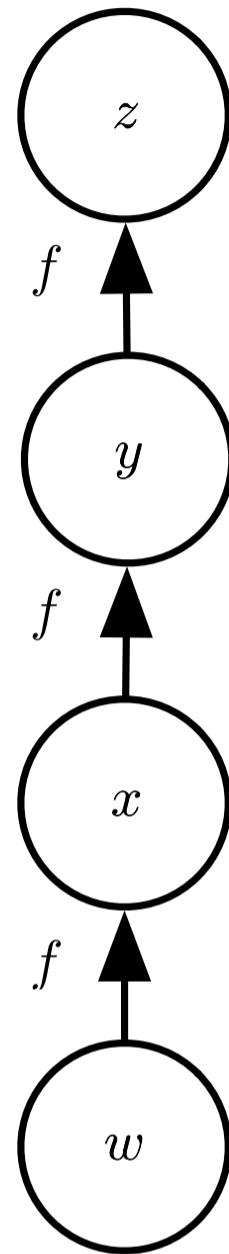
- Suppose that  $\mathbf{y} = g(\mathbf{x})$  and  $z = f(\mathbf{y})$ , then  $\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$
- In vector notation, this can be equivalently rewritten as:

$$\nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

Back-propagation consists of performing one Jacobian-gradient product for each operation in the graph

# Back-propagation as Recursive Chain Rule

- Using the chain rule, it is straightforward to write down an expression for the gradient of a scalar w.r.t. **any node in the graph**
- Actually evaluating that expression involves some extra considerations
- Specifically, many subexpressions may be repeated several times within overall expression for gradient



$$\begin{aligned}\frac{\partial z}{\partial w} &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \\ &= f'(y) f'(x) f'(w) \\ &= f'(f(f(w))) f'(f(w)) f'(w)\end{aligned}$$

# Back-propagation in a Computational Graph

- Consider a computational graph describing how to compute a **single scalar**  $u^{(n)}$  (e.g. the loss on a single example)
- This scalar is the quantity whose gradient we want to obtain, w.r.t. the  $n_i$  input nodes,  $\{u^{(1)}, \dots, u^{(n_i)}\}$  (e.g. the parameters of a model)
  - In other words, we wish to compute  $\frac{\partial u^{(n)}}{\partial u^{(i)}} \forall i \in \{1, 2, \dots, n_i\}$
- Assume that the nodes of the graph are ordered in such a way that we can compute their output one after another, starting at  $u^{(n_i+1)}$  and going up to  $u^{(n)}$
- Each node  $u^{(i)}$  is associated with an operation  $f^{(i)}$  and is computed by evaluating:

$$u^{(i)} = f^{(i)}(\mathbb{A}^{(i)})$$

set of all nodes that are parents of  $u^{(i)}$

- Back-propagation computes each derivative  $\frac{\partial u^{(n)}}{\partial u^{(i)}}$  associated with the respective forward graph node  $u^{(i)}$ , via the chain rule:

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i:j \in Pa(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \frac{\partial u^{(i)}}{\partial u^{(j)}}$$

# Simplified Forward Pass

(See Algorithm 6.1 in Goodfellow et al.)

# Simplified Backward Pass

(See Algorithm 6.2 in Goodfellow et al.)

# Forward Propagation in MLP

Example of a 2 hidden layer (or 4 layer) network ...



|           |                                    |
|-----------|------------------------------------|
| $x$       | input                              |
| $h^{(1)}$ | 1 <sup>st</sup> layer hidden units |
| $h^{(2)}$ | 2 <sup>nd</sup> layer hidden units |
| $\hat{y}$ | output (prediction)                |

# Forward Propagation (L1)

Forward propagation is the process of computing the output of the network given its input



$$\mathbf{x} \in \mathbb{R}^d \quad \mathbf{W}^{(1)} \in \mathbb{R}^{n_1 \times d} \quad \mathbf{b}^{(1)} \in \mathbb{R}^{n_1} \quad \mathbf{h}^{(1)} \in \mathbb{R}^{n_1}$$

$$\mathbf{a}^{(1)} = \mathbf{b}^{(1)} + \mathbf{W}^{(1)}\mathbf{x}$$

$$\mathbf{h}^{(1)} = g(\mathbf{a}^{(1)})$$

$$\mathbf{W}^{(1)}$$

1<sup>st</sup> layer weights

$$g(\cdot)$$

activation function  
(e.g. ReLU)

$$\mathbf{b}^{(1)}$$

1<sup>st</sup> layer biases

# Forward Propagation (L2)



$$h^{(1)} \in \mathbb{R}^{n_1} \quad W^{(2)} \in \mathbb{R}^{n_2 \times n_1} \quad b^{(2)} \in \mathbb{R}^{n_2} \quad h^{(2)} \in \mathbb{R}^{n_2}$$

$$a^{(2)} = b^{(2)} + W^{(2)}h^{(1)}$$

$$h^{(2)} = g(a^{(2)})$$

$W^{(2)}$       2<sup>nd</sup> layer weights

$b^{(2)}$       2<sup>nd</sup> layer biases

# Forward Propagation (Out)



$$h^{(2)} \in \mathbb{R}^{n_2} \quad W^{(3)} \in \mathbb{R}^{n_3 \times n_2} \quad b^{(3)} \in \mathbb{R}^{n_3} \quad \hat{y} \in \mathbb{R}^{n_3}$$

$$\mathbf{a}^{(3)} = \mathbf{b}^{(3)} + \mathbf{W}^{(3)} \mathbf{h}^{(2)}$$

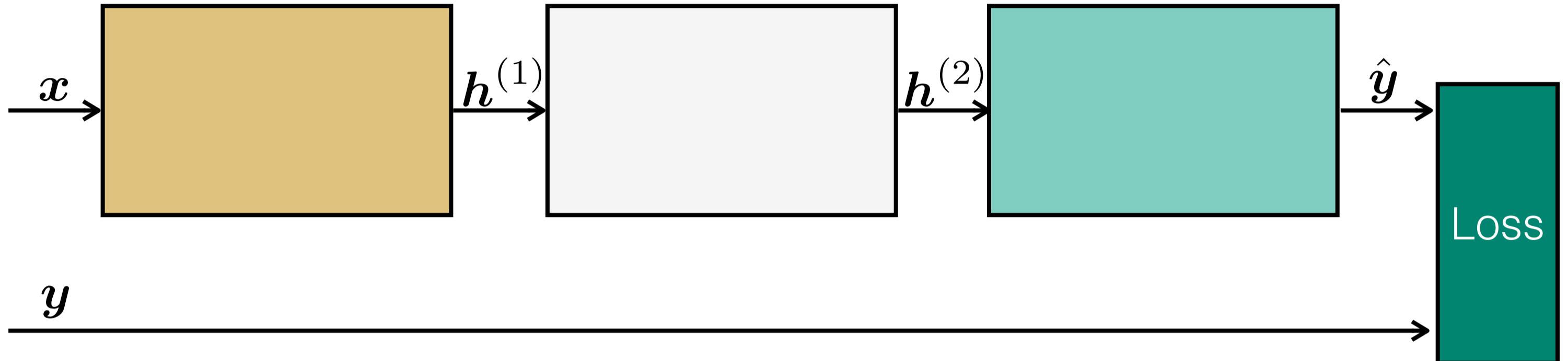
$$\hat{y} = g(\mathbf{a}^{(3)})$$

$W^{(3)}$       3<sup>rd</sup> layer weights

$b^{(3)}$       3<sup>rd</sup> layer biases

Note:  $g(\cdot)$   
is usually different  
at the output layer

# How Good is a Network?



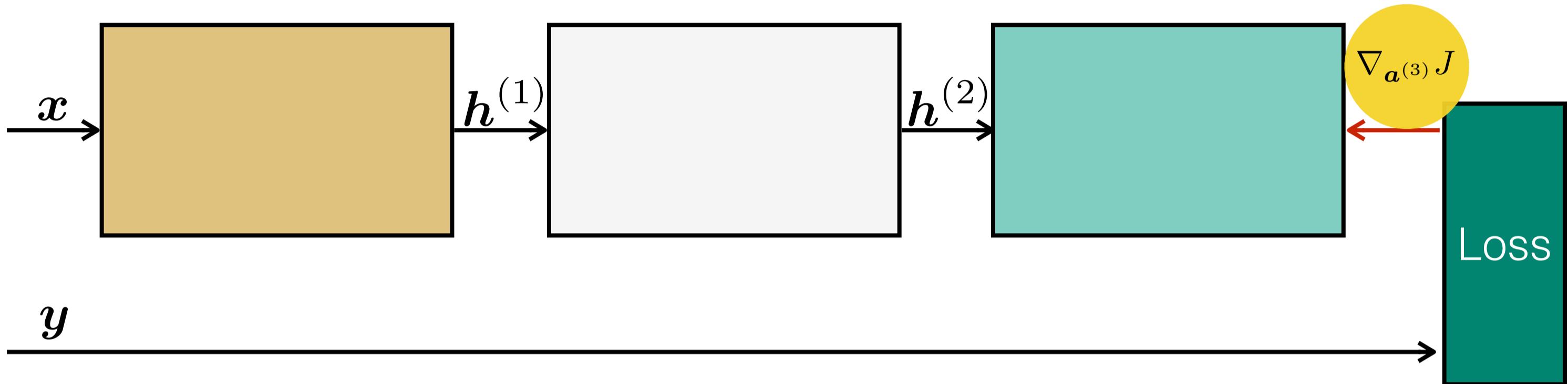
$y$  target (ground truth)

$\theta \equiv \{W^{(1)}, b^{(1)}, \dots, W^{(l)}, b^{(l)}\}$  parameters

$L(\hat{y}, y)$  loss function

$J = L(\hat{y}, y) + \alpha \Omega(\theta)$  regularized cost

# Backward Propagation (Out)



After the forward computation, compute the gradient on the output layer

$$g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$$

Then, convert the gradient on the layer's output into a gradient on the pre-activation:

$$g \leftarrow \nabla_{a^{(3)}} J$$

# Derivative w.r.t. Input of Softmax

First convert the discrepancy between each output and its target value into an error derivative

$$\hat{y}_k = p(c_k = 1 | \mathbf{x}) = \frac{e^{a_k}}{\sum_{j=1}^C e^{a_j}}$$

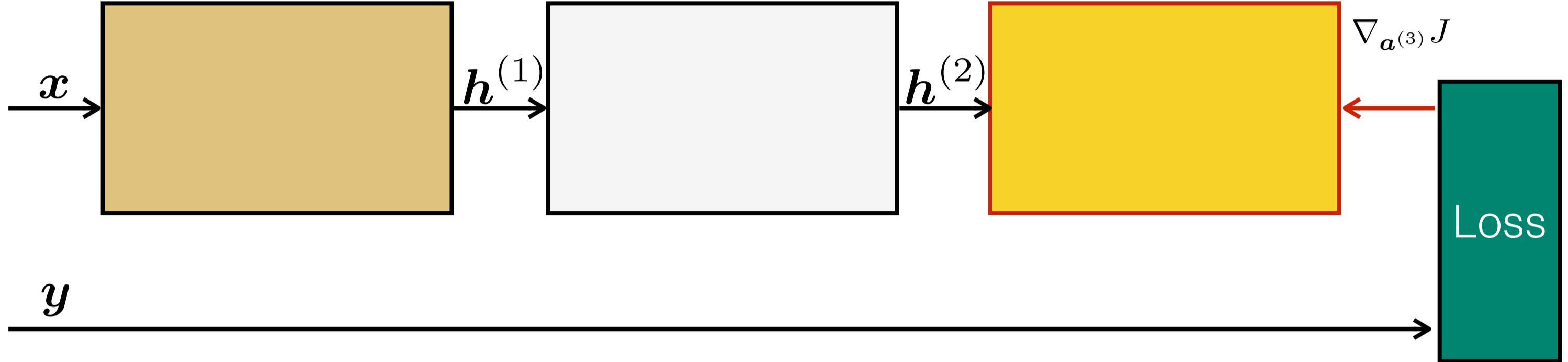
Note: layer superscript has been left off for simplicity, i.e.  $a_j = (\mathbf{a}^{(3)})_j$

$$L(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_j y_j \log p(c_j | \mathbf{x}) \quad \mathbf{y} = [00 \dots 010 \dots 0]$$

By substituting the first formula into the second, and taking the derivative w.r.t.  $a$  we get:

$$\frac{\partial L}{\partial a_j} = p(c_j | \mathbf{x}) - y_j$$

# Backward Propagation (Update)

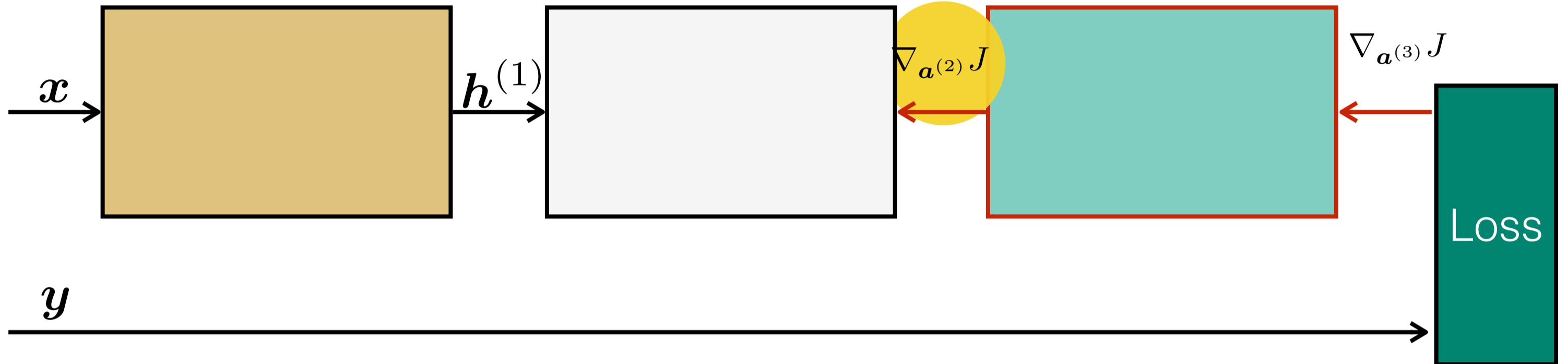


Compute gradients on weights and biases:

$$\nabla_{\mathbf{b}^{(3)}} J = \mathbf{g} + \alpha \nabla_{\mathbf{b}^{(3)}} \Omega(\boldsymbol{\theta})$$

$$\nabla_{\mathbf{W}^{(3)}} J = \mathbf{g} \mathbf{h}^{(2)\top} + \alpha \nabla_{\mathbf{W}^{(3)}} \Omega(\boldsymbol{\theta})$$

# Backward Propagation (L2)



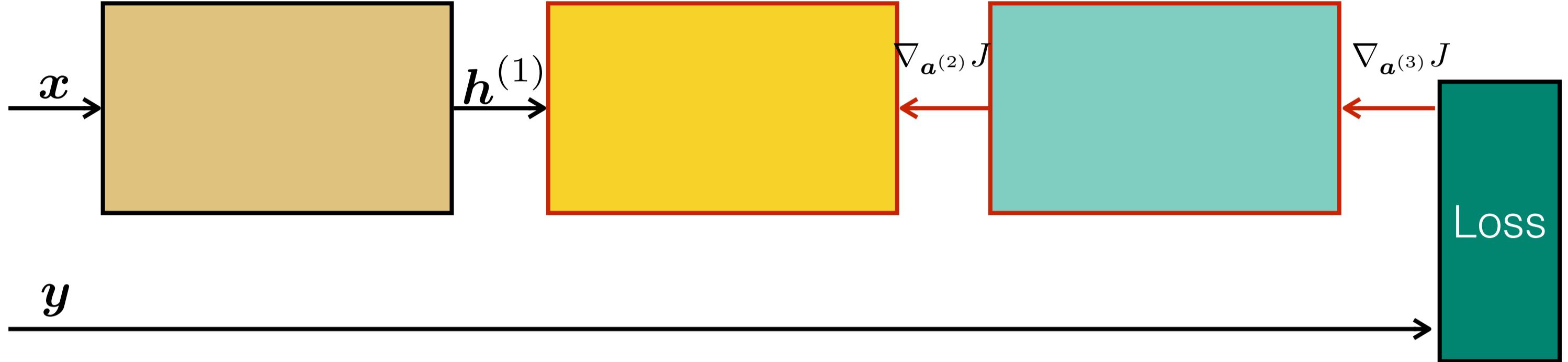
Propagate the gradients w.r.t. the next lower-level hidden layer's activation

$$g \leftarrow \nabla_{h^{(2)}} J = W^{(3)\top} g$$

Then, convert the gradient on the layer's output into a gradient on the pre-activation (element-wise nonlinearity):

$$g \leftarrow \nabla_{\mathbf{a}^{(2)}} J = g \odot g'(\mathbf{a}^{(2)})$$

# Backward Propagation (Update)

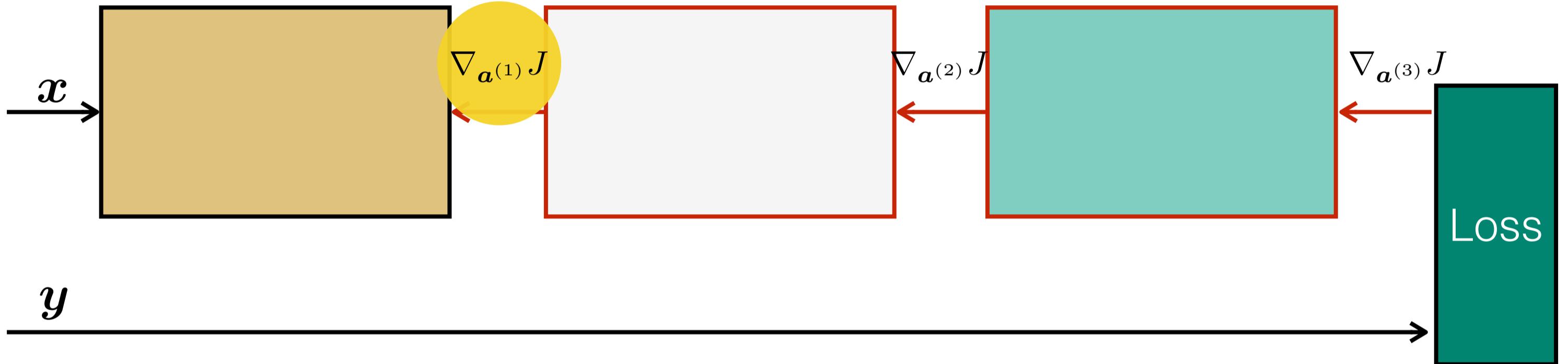


Compute gradients on weights and biases:

$$\nabla_{\mathbf{b}^{(2)}} J = \mathbf{g} + \alpha \nabla_{\mathbf{b}^{(2)}} \Omega(\boldsymbol{\theta})$$

$$\nabla_{\mathbf{W}^{(2)}} J = \mathbf{g} \mathbf{h}^{(1)\top} + \alpha \nabla_{\mathbf{W}^{(2)}} \Omega(\boldsymbol{\theta})$$

# Backward Propagation (L1)



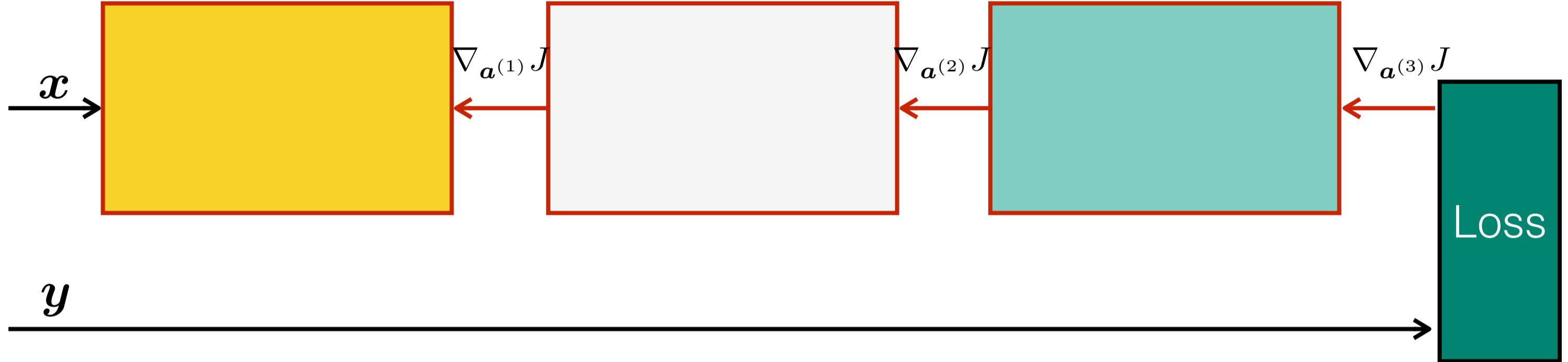
Propagate the gradients w.r.t. the next lower-level hidden layer's activation

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(1)}} J = \mathbf{W}^{(2)\top} \mathbf{g}$$

Then, convert the gradient on the layer's output into a gradient on the pre-activation (element-wise nonlinearity):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(1)}} J = \mathbf{g} \odot g'(\mathbf{a}^{(1)})$$

# Backward Propagation (Update)



Compute gradients on weights and biases:

$$\nabla_{\mathbf{b}^{(1)} J} = \mathbf{g} + \alpha \nabla_{\mathbf{b}^{(1)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(1)} J} = \mathbf{g} \mathbf{x}^\top + \alpha \nabla_{\mathbf{W}^{(1)}} \Omega(\theta)$$