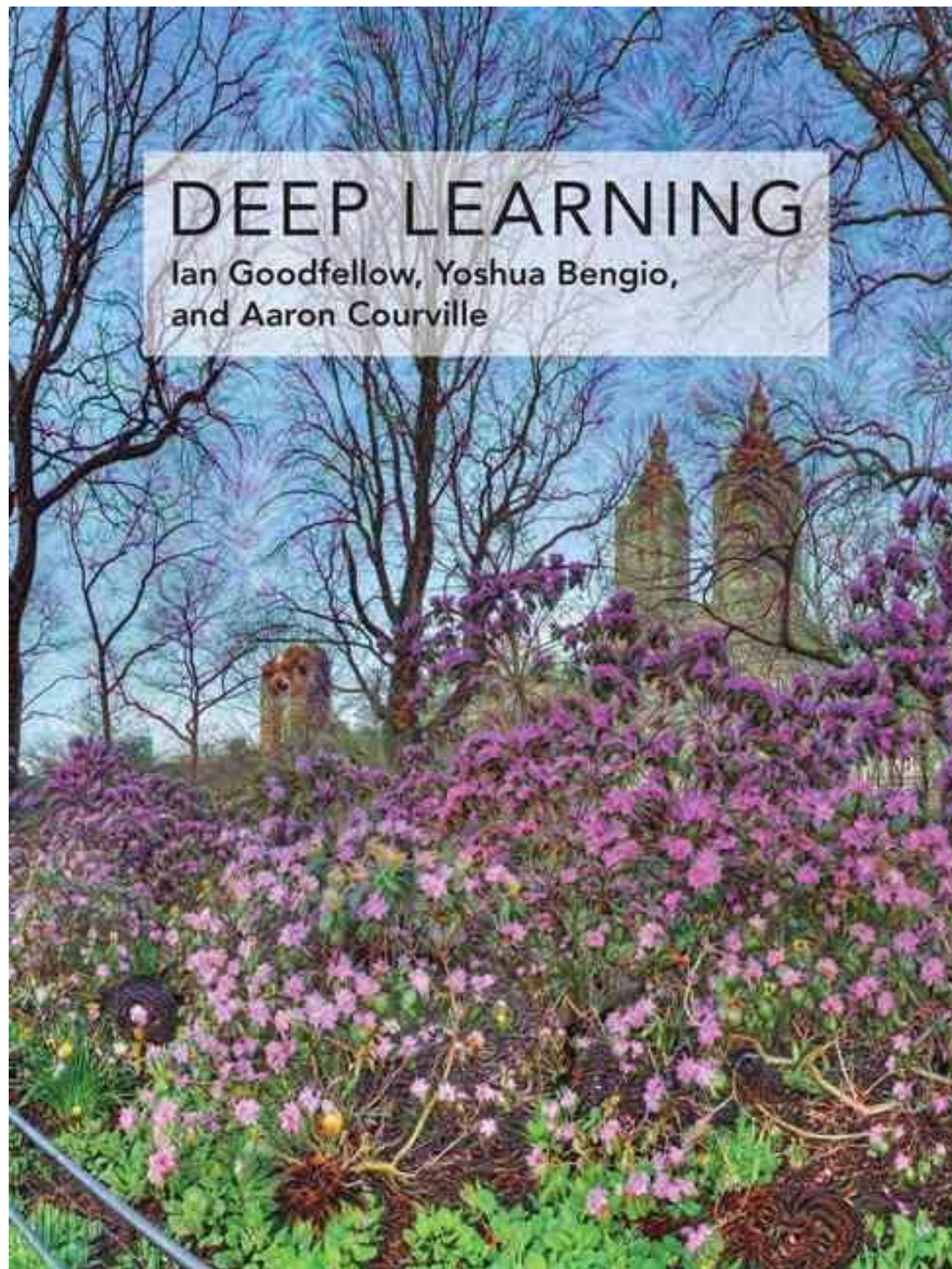


Regularization, SGD, and Backpropagation

Bill Freeman, Antonio Torralba, Phillip Isola
6.819 / 6.869



<http://www.deeplearningbook.org/>

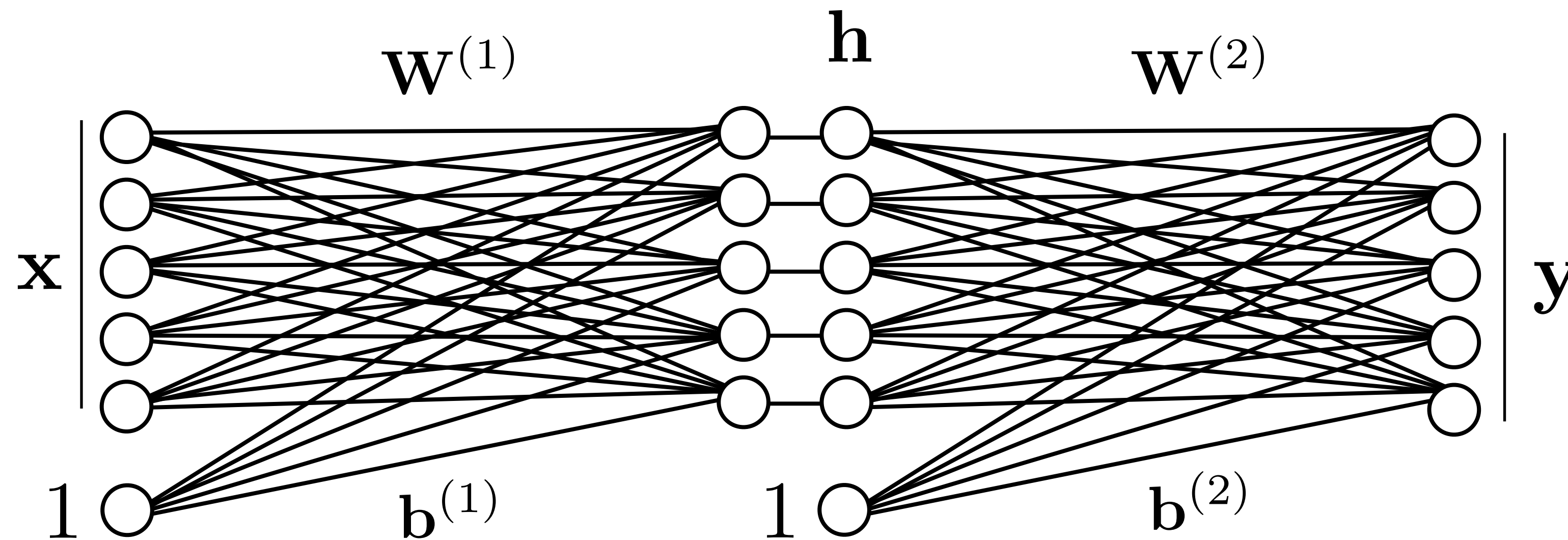
By Ian Goodfellow, Yoshua Bengio and Aaron Courville

November 2016

Today: parts of chapters 7 and 8

Stacking layers

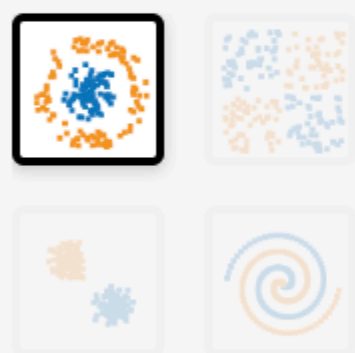
Input representation Intermediate representation Output representation



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

DATA

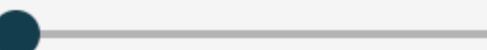
Which dataset do you want to use?



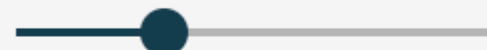
Ratio of training to test data: 50%



Noise: 0



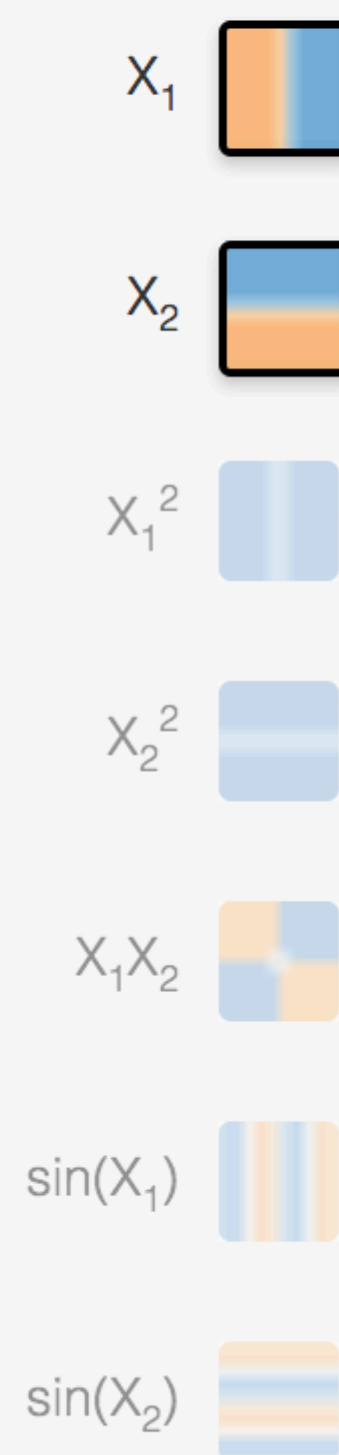
Batch size: 10



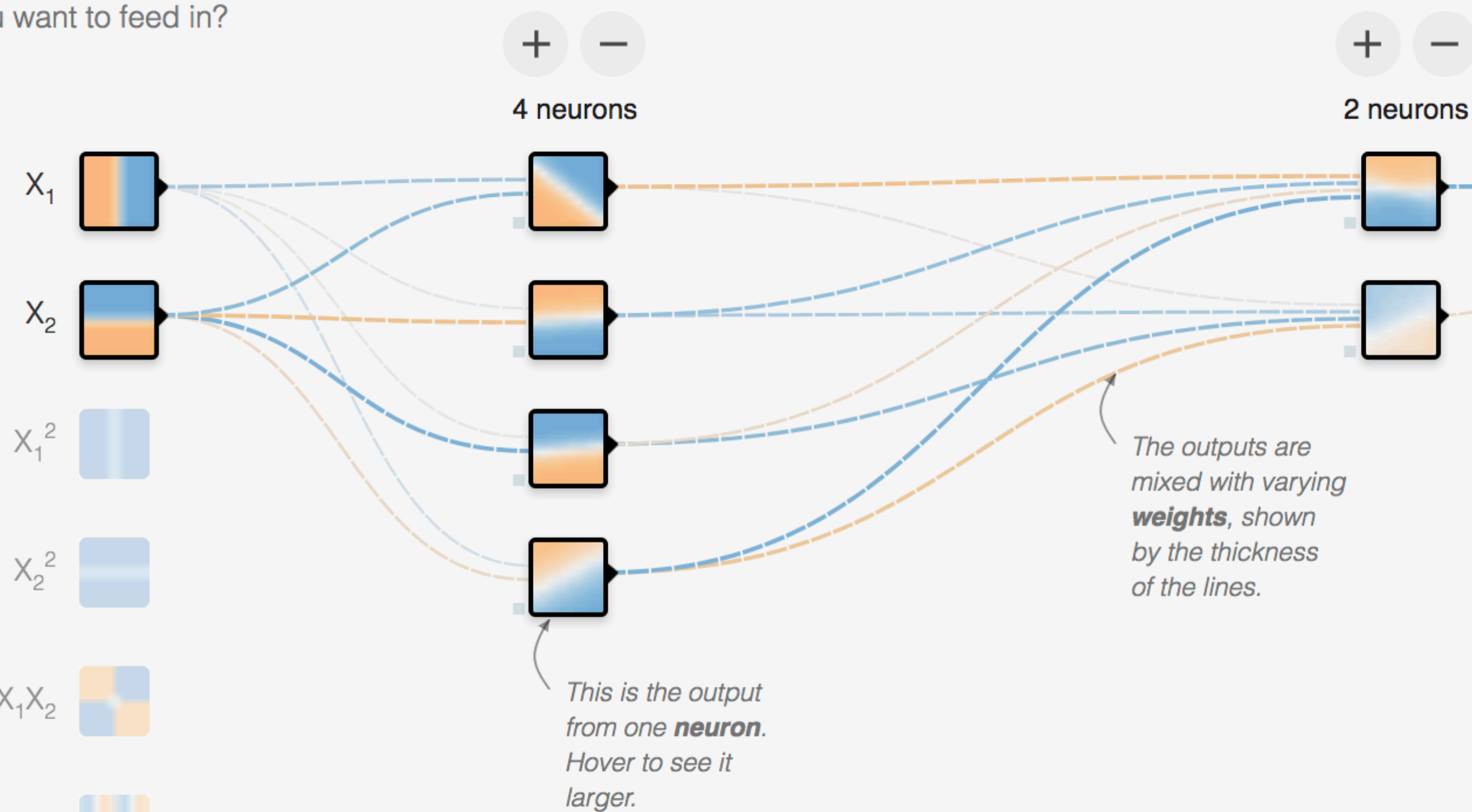
REGENERATE

FEATURES

Which properties do you want to feed in?



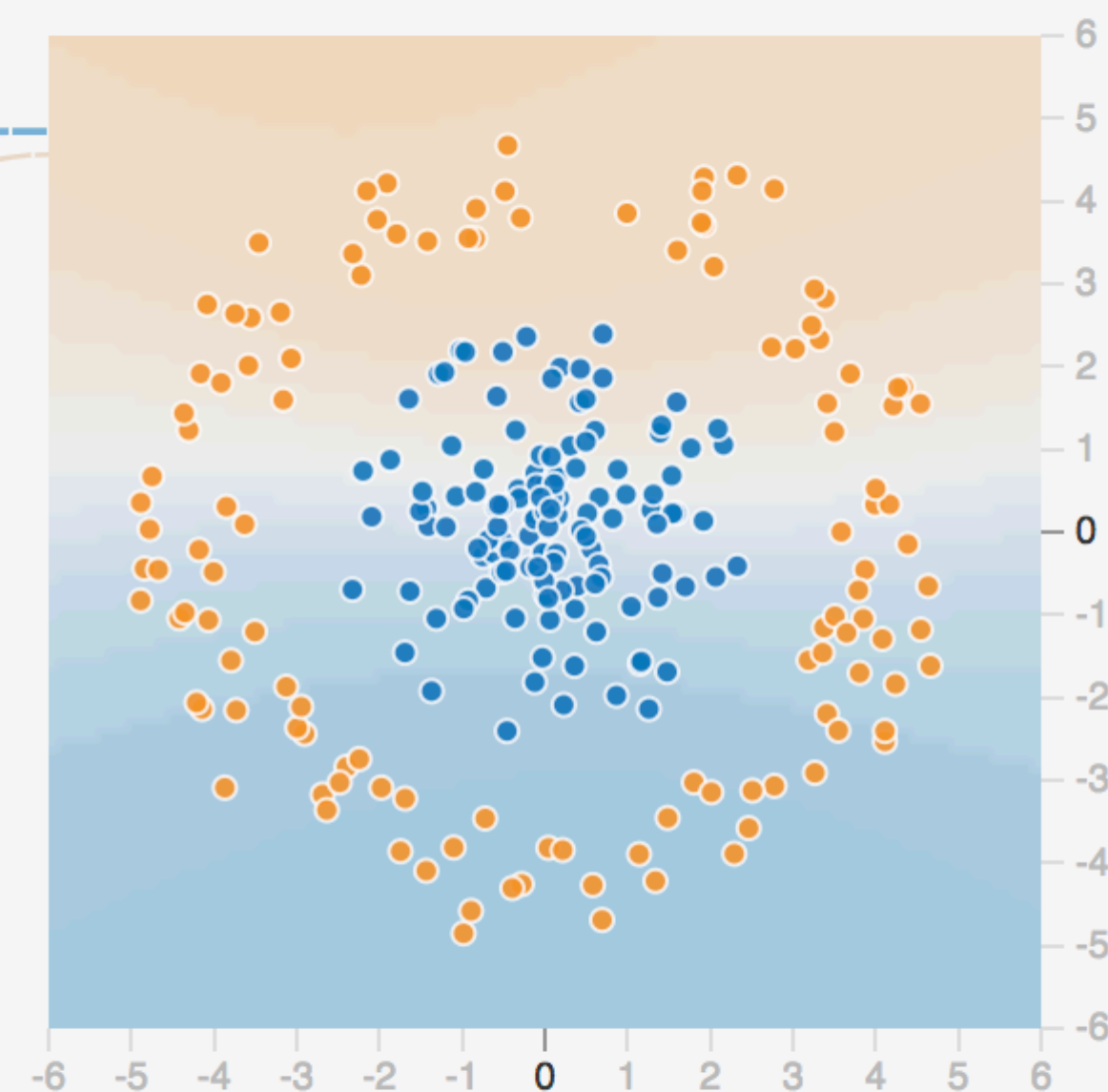
+ - 2 HIDDEN LAYERS



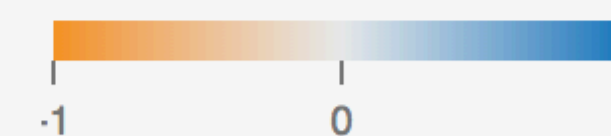
OUTPUT

Test loss 0.540

Training loss 0.555



Colors shows data, neuron and weight values.

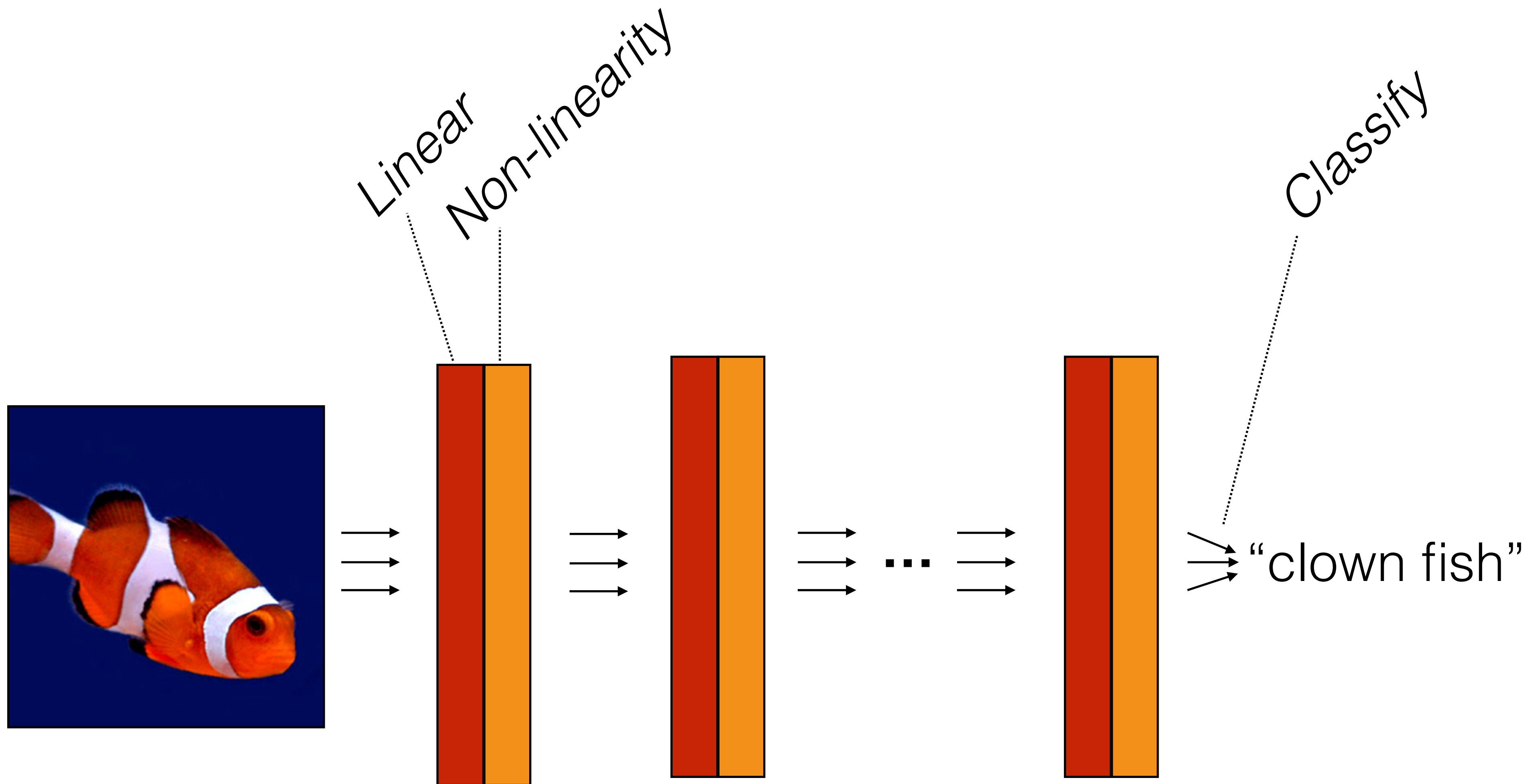


☐ Show test data

☐ Discretize output

[<http://playground.tensorflow.org>]

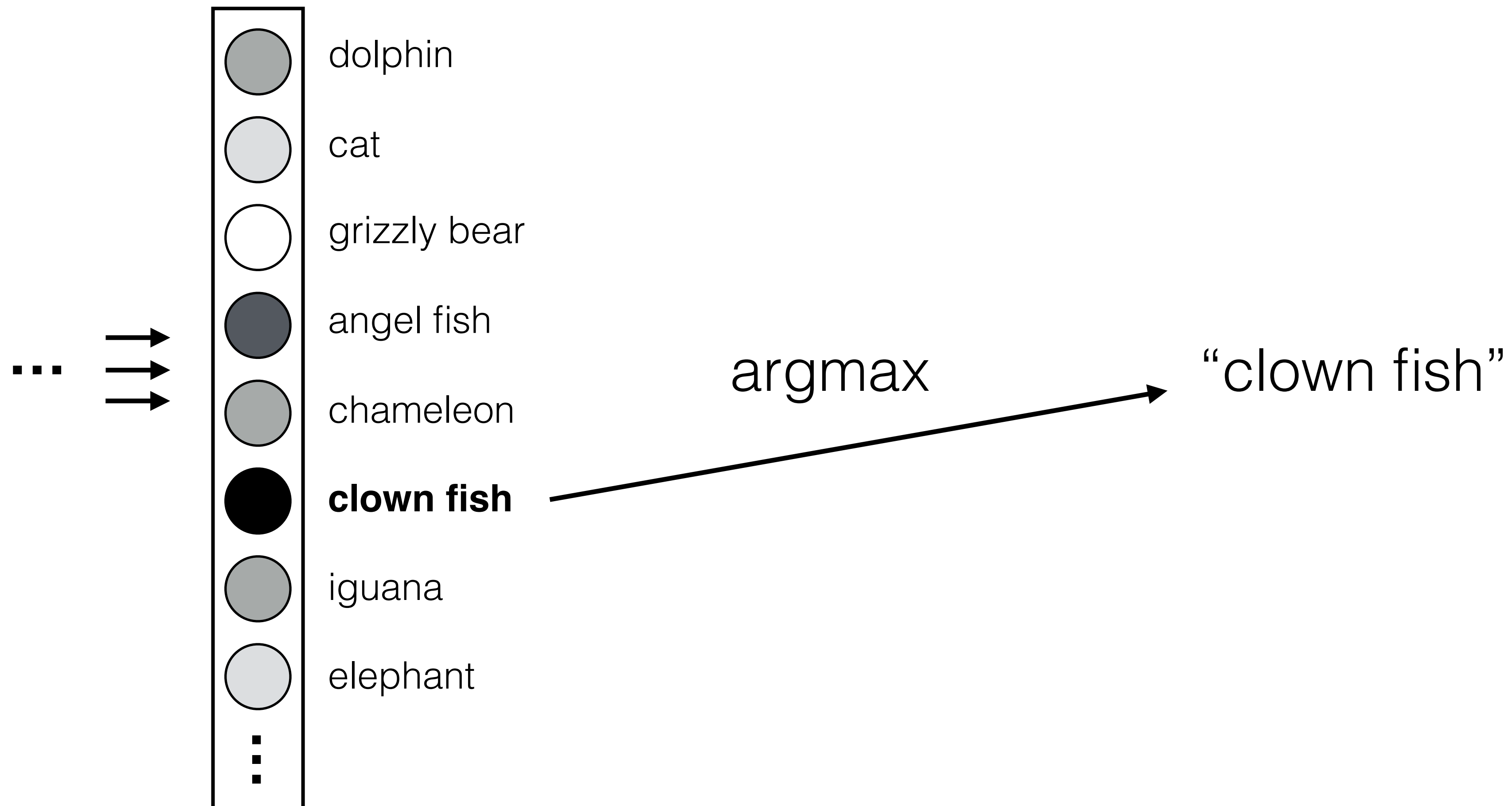
Deep nets



$$f(\mathbf{x}) = f^{(L)}(\dots f^{(2)}(f^{(1)}(\mathbf{x})))$$

Classifier layer

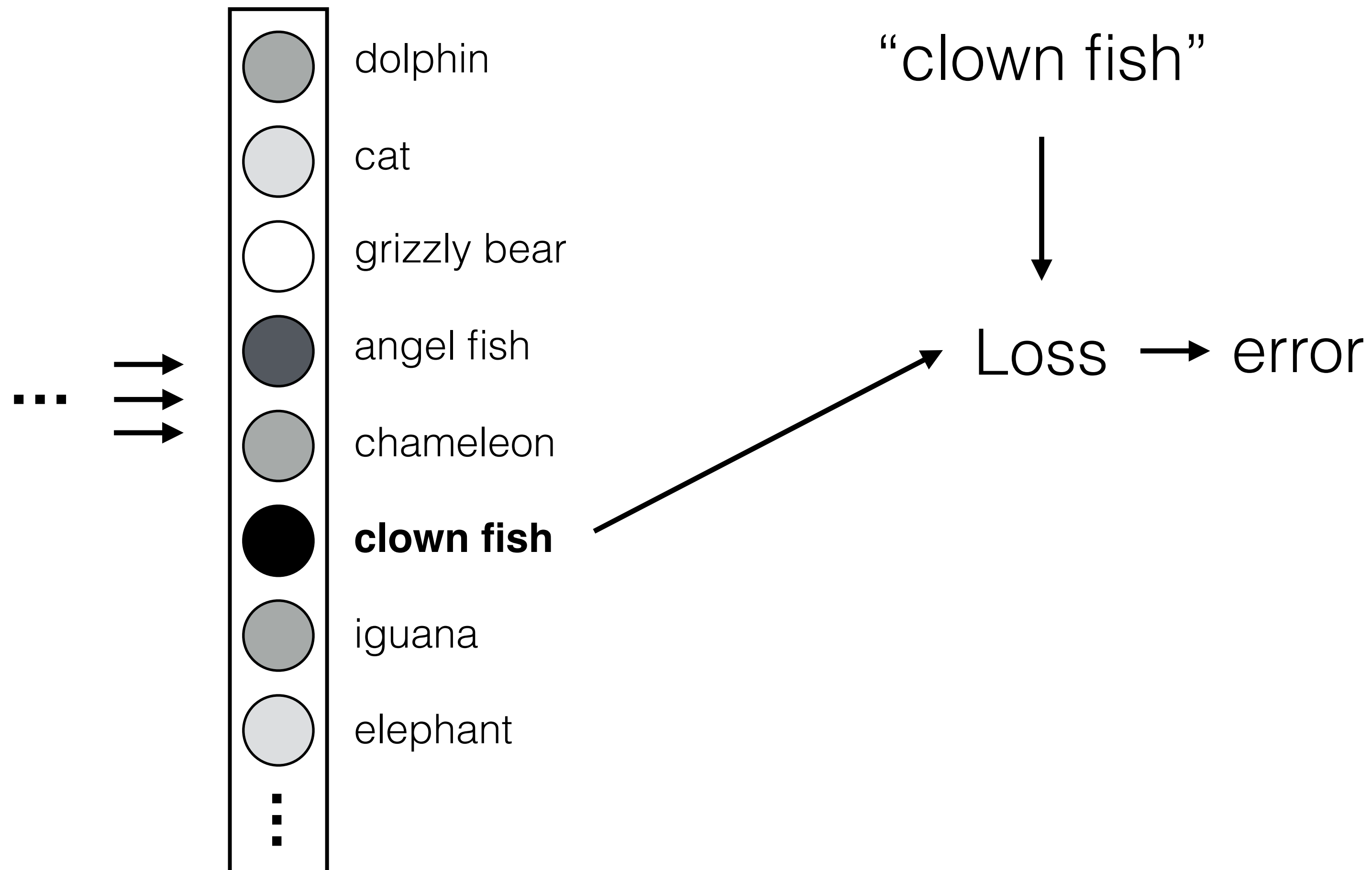
Last layer



Loss function

Network output

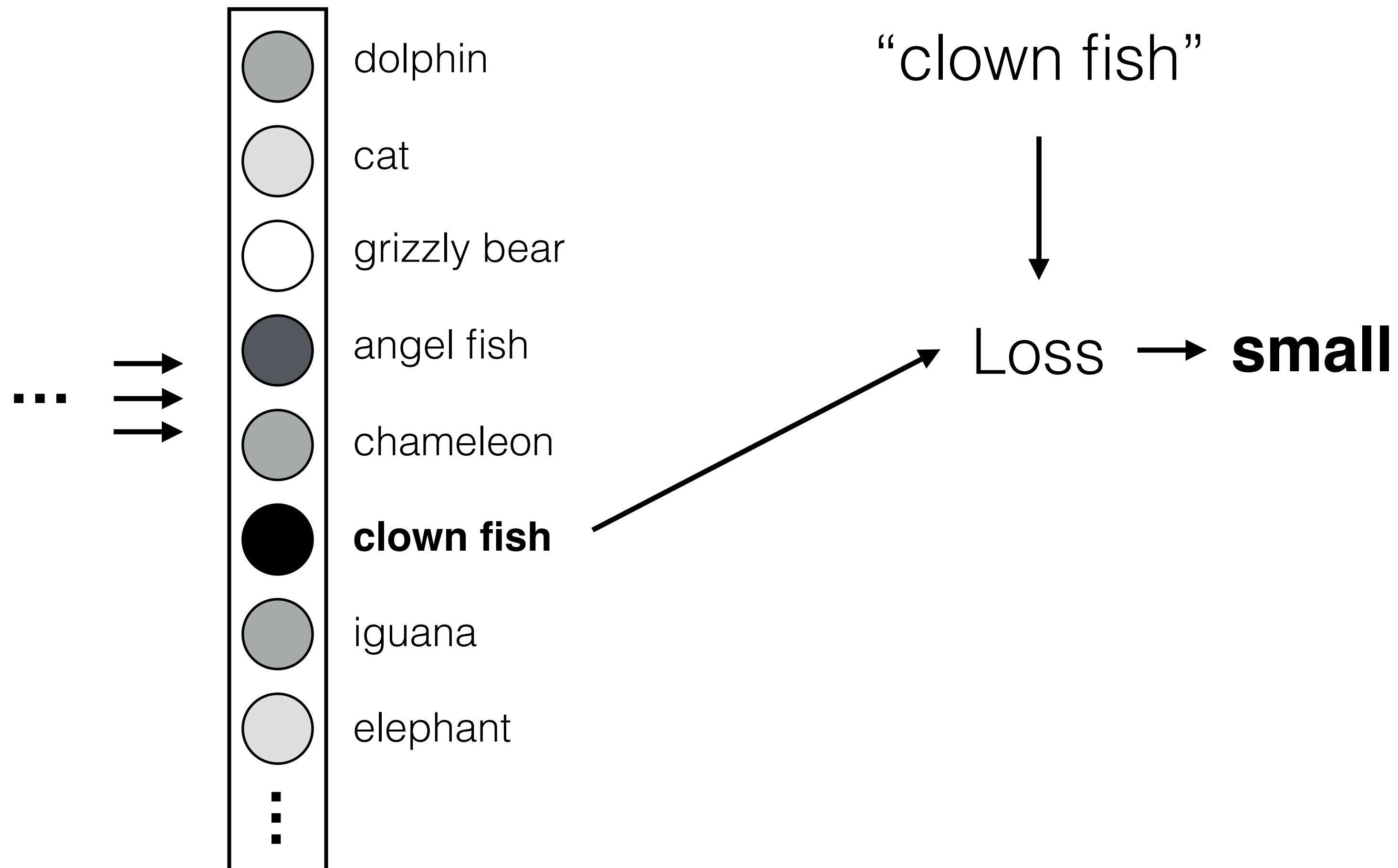
Ground truth label



Loss function

Network output

Ground truth label

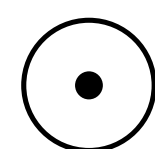
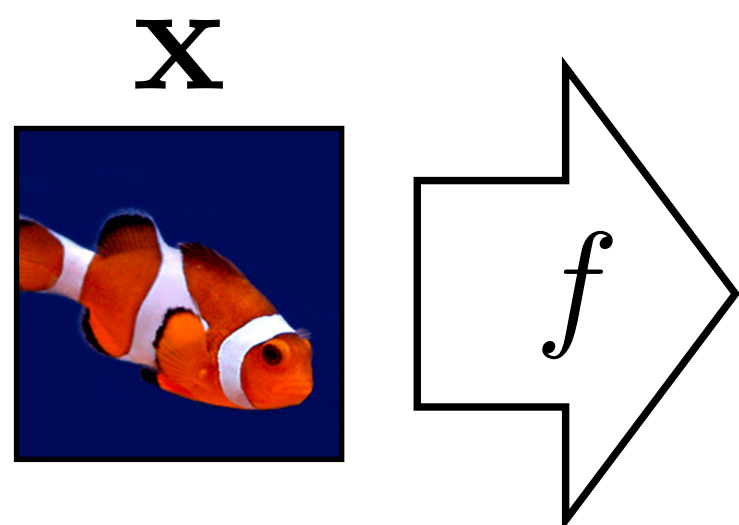


Loss function

Network output

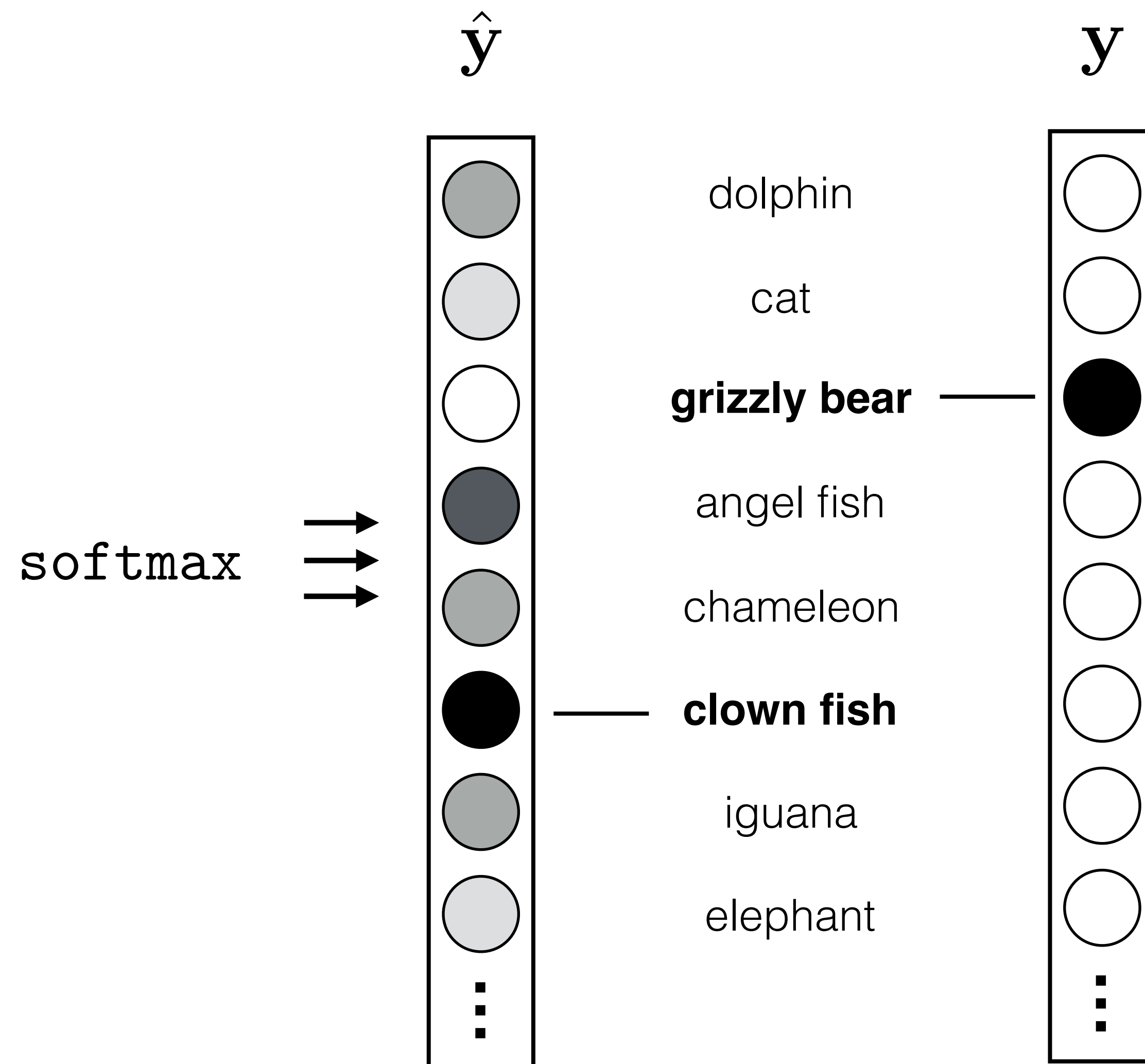
Ground truth label





Network output

Ground truth label



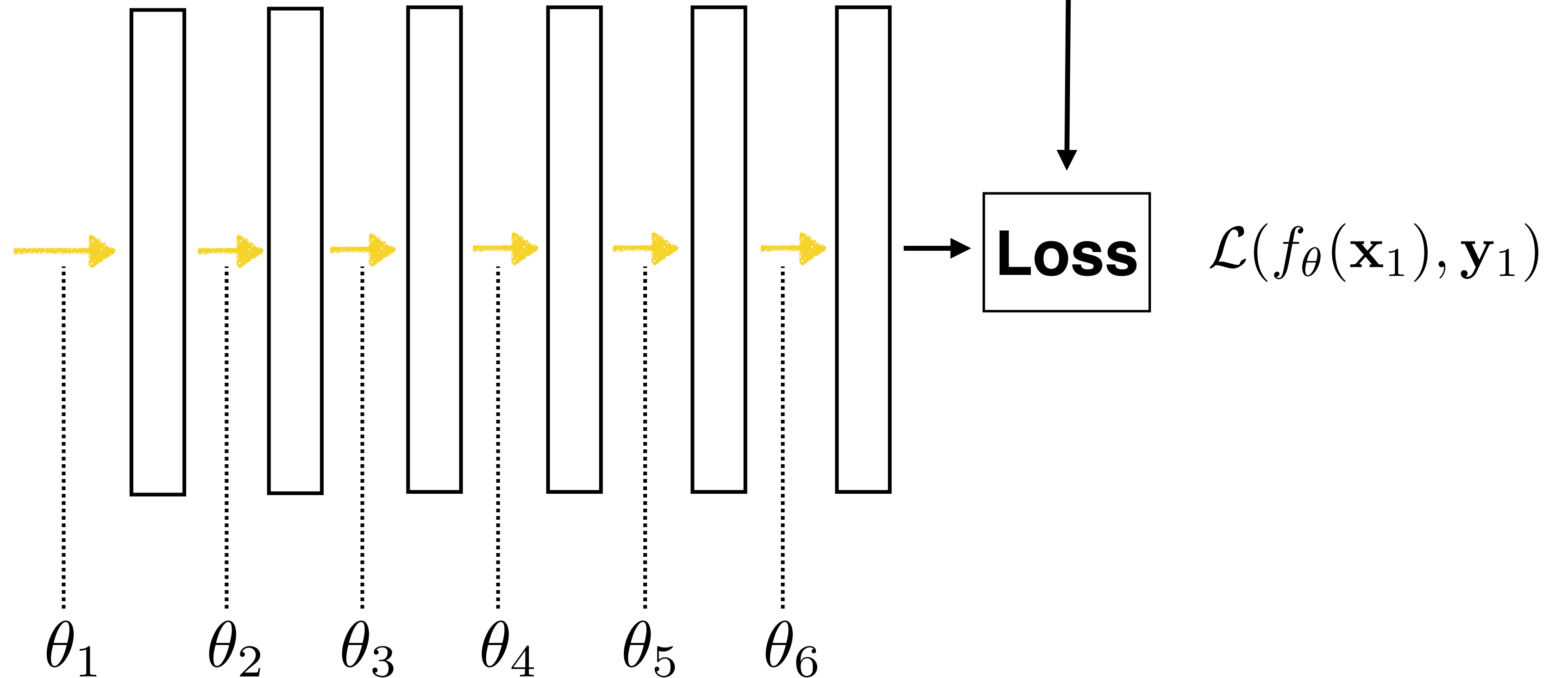
Probability of the observed
data under the model

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$

Deep learning

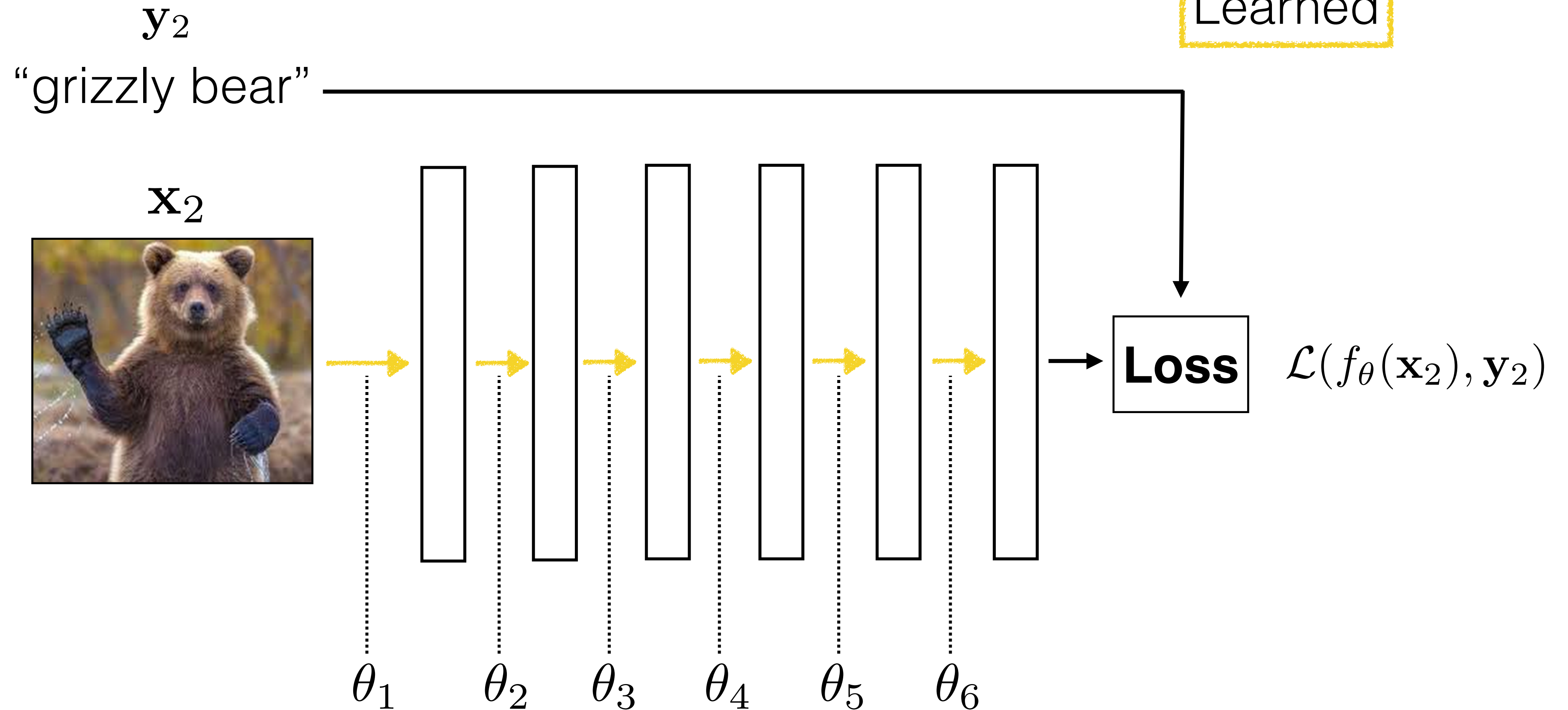
\mathbf{y}_1
“clown fish”

Learned



$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

Deep learning

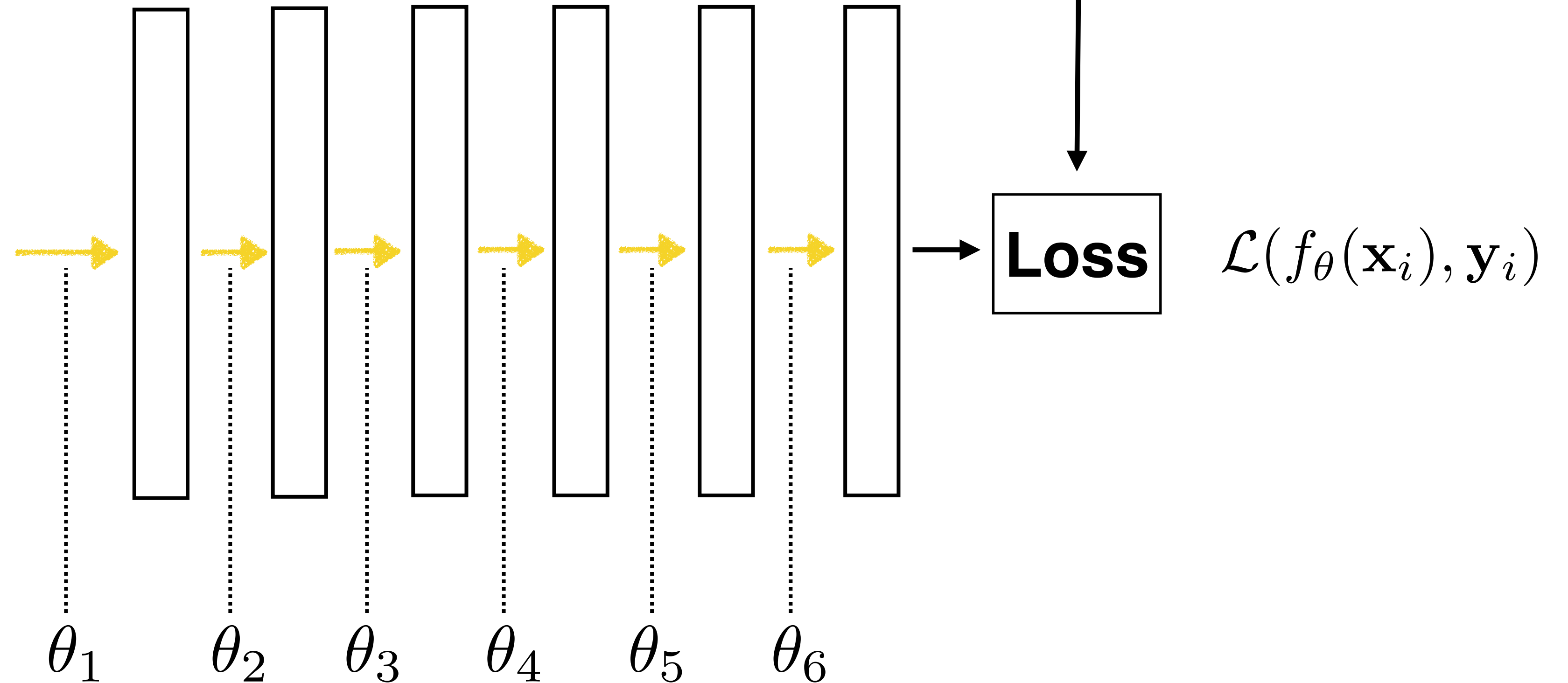


$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

Deep learning

y_i
“chameleon”

Learned



$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i)$$

Regularizing deep nets

Deep nets have millions of parameters!

On many datasets, it is easy to overfit — we may have more free parameters than data points to constrain them.

How can we regularize to prevent the network from overfitting?

1. Fewer neurons, fewer layers
2. Weight decay
3. Dropout
4. Normalization layers
5. ...

Recall: regularized least squares

$$f_{\theta}(x) = \sum_{k=0}^K \theta_k x^k$$

$$R(\theta) = \lambda \|\theta\|_2^2 \longleftarrow \text{Only use polynomial terms if you really need them! Most terms should be zero}$$

ridge regression, a.k.a., **Tikhonov regularization**

Probabilistic interpretation: R is a Gaussian **prior** over values of the parameters.

Regularizing the weights in a neural net

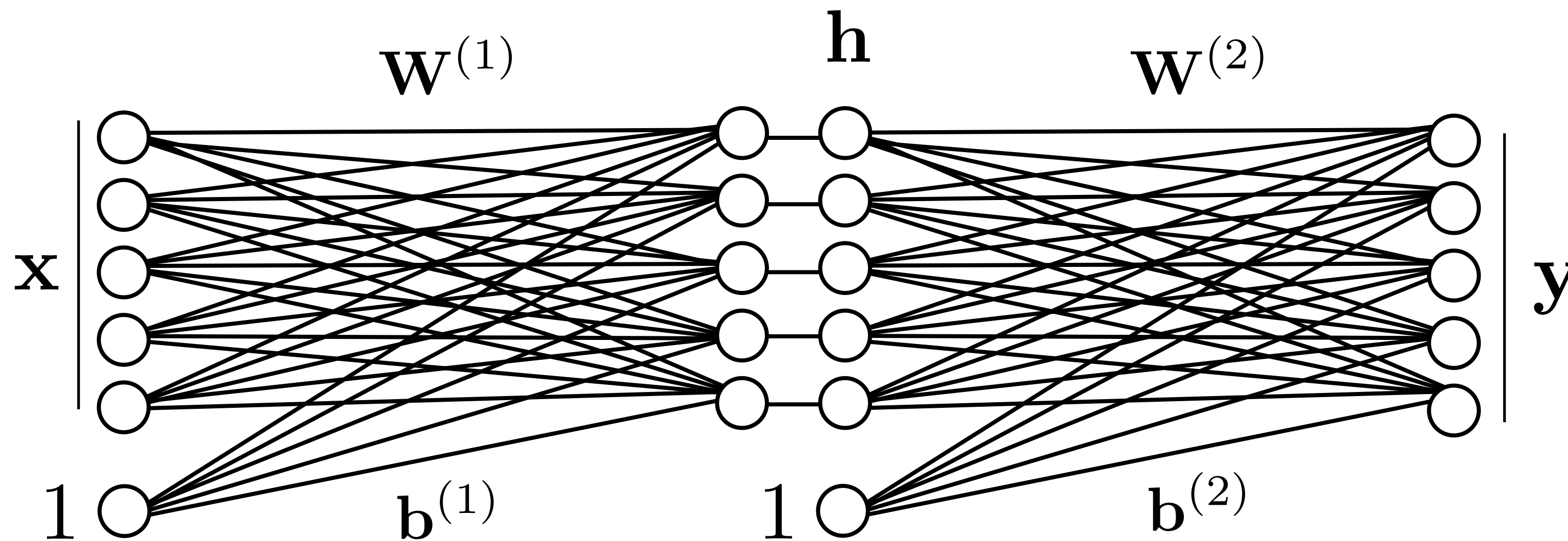
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + R(\theta)$$

$$R(\mathbf{W}) = \lambda \|\mathbf{W}\|_2^2 \quad \longleftarrow \quad \textbf{weight decay}$$

“We prefer to keep weights small.”

Dropout

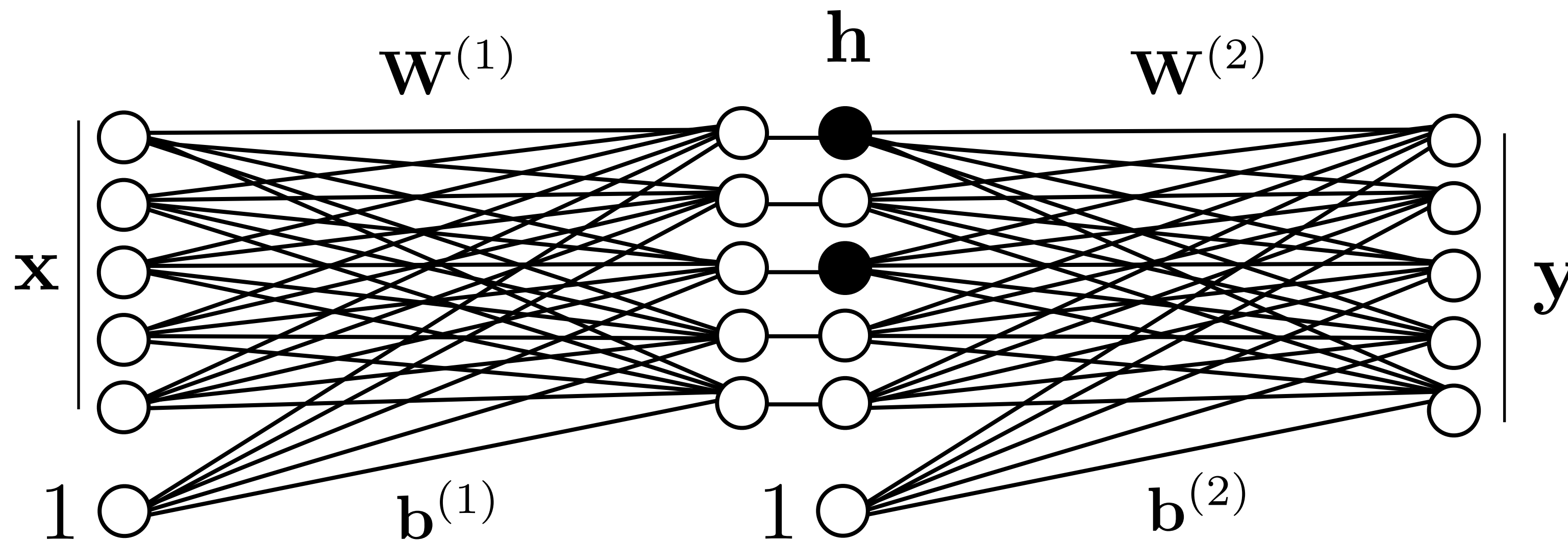
Input representation Intermediate representation Output representation



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

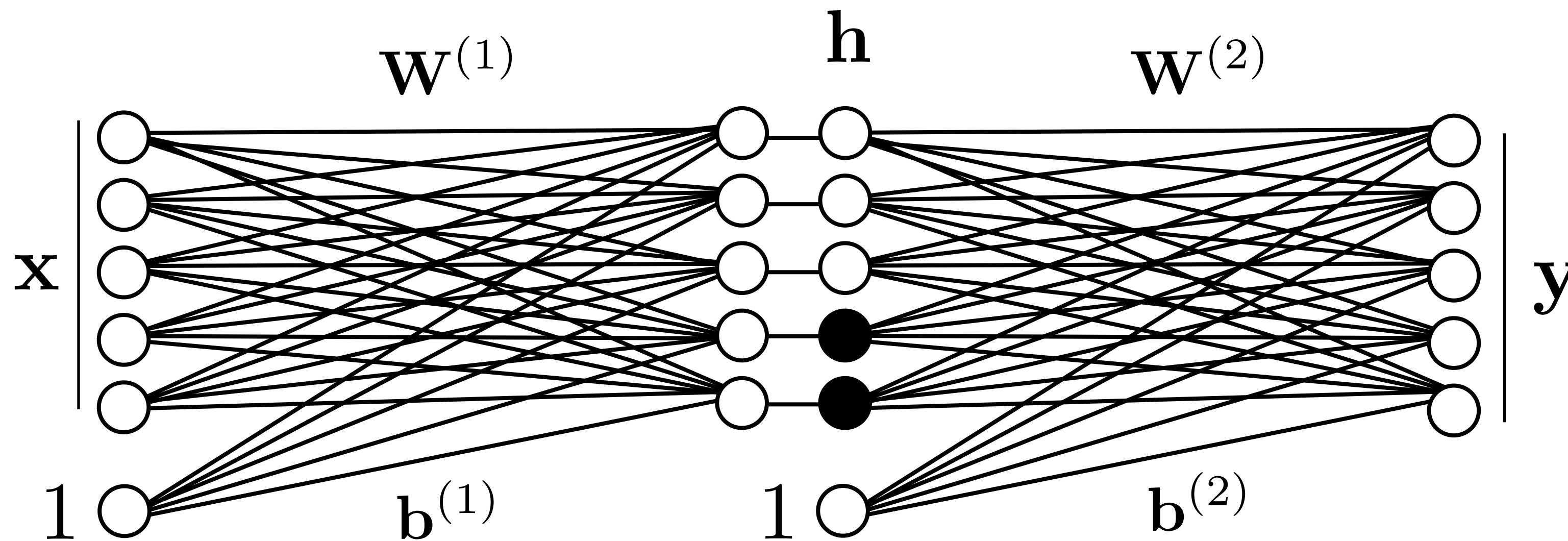
Input representation Intermediate representation Output representation



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

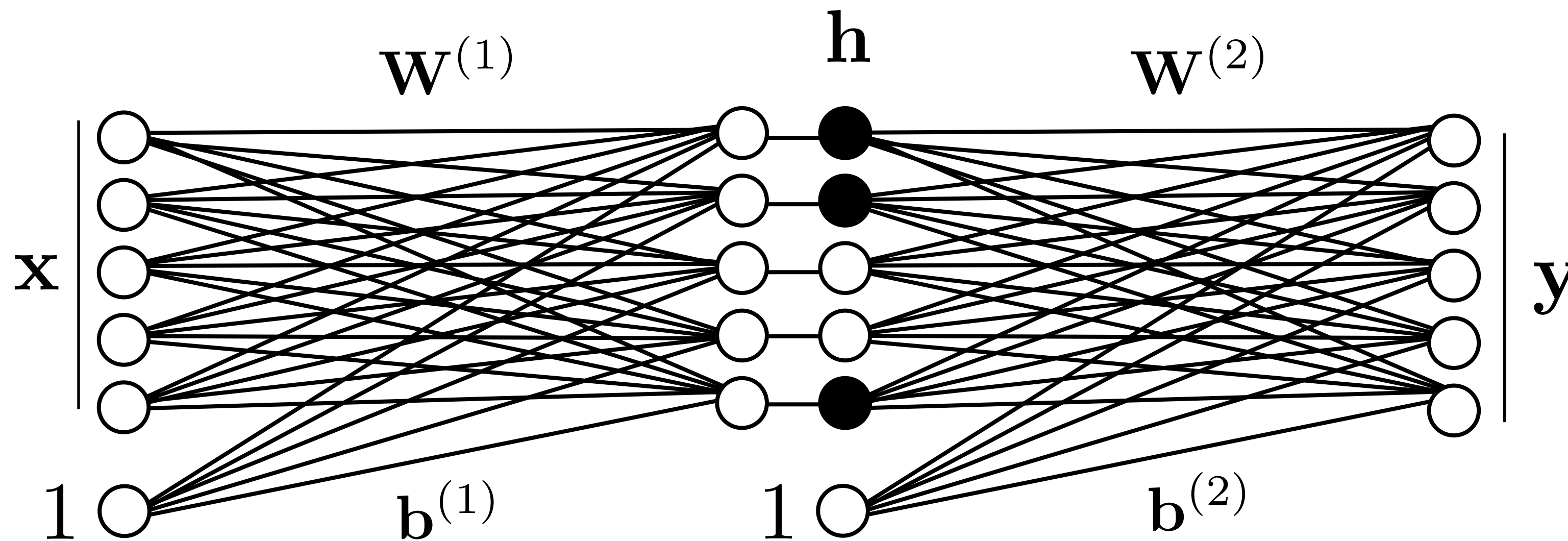
Input representation Intermediate representation Output representation



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

Dropout

Input representation Intermediate representation Output representation



$$\theta = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L)}\}$$

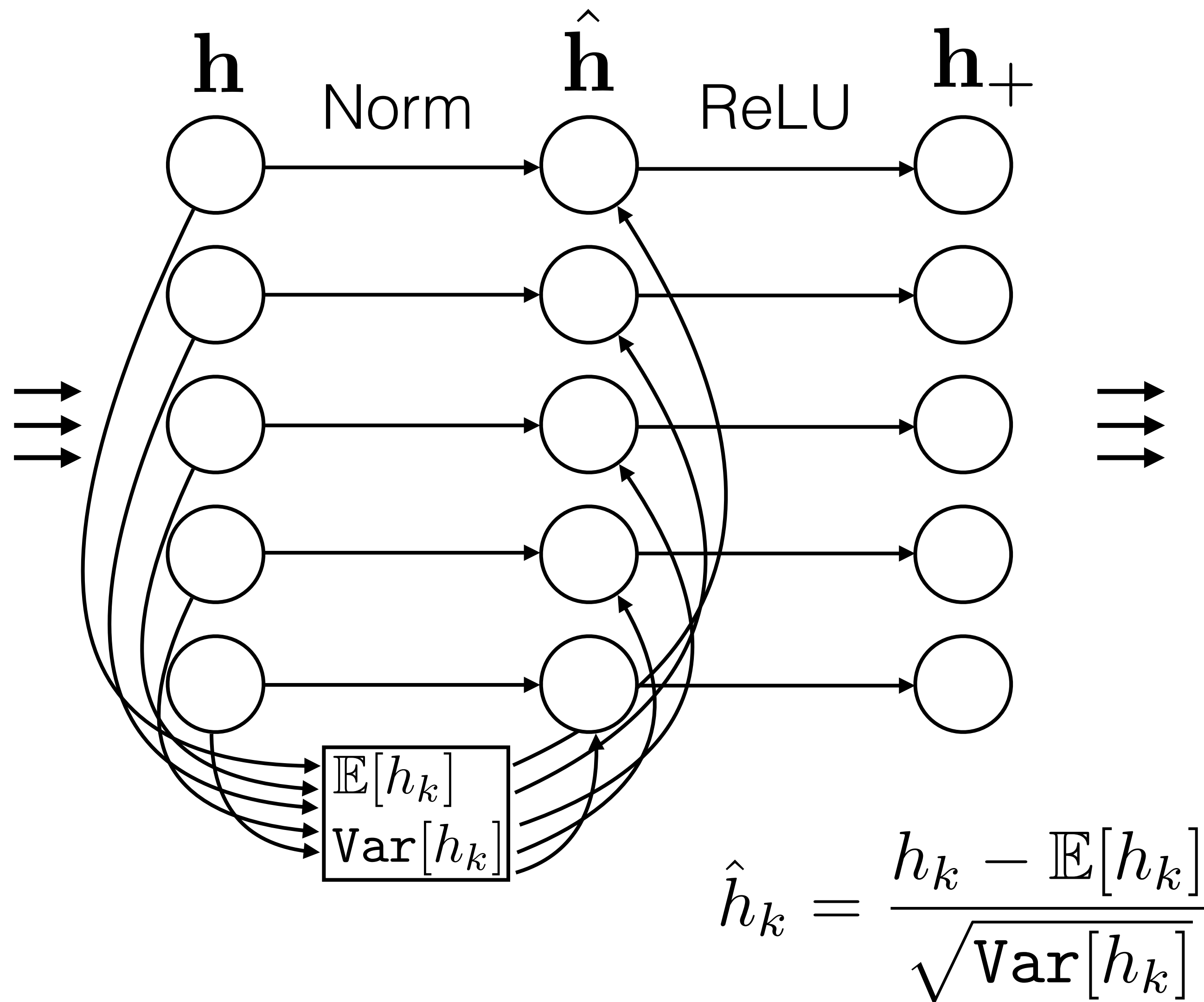
Dropout

Randomly zero out hidden units.

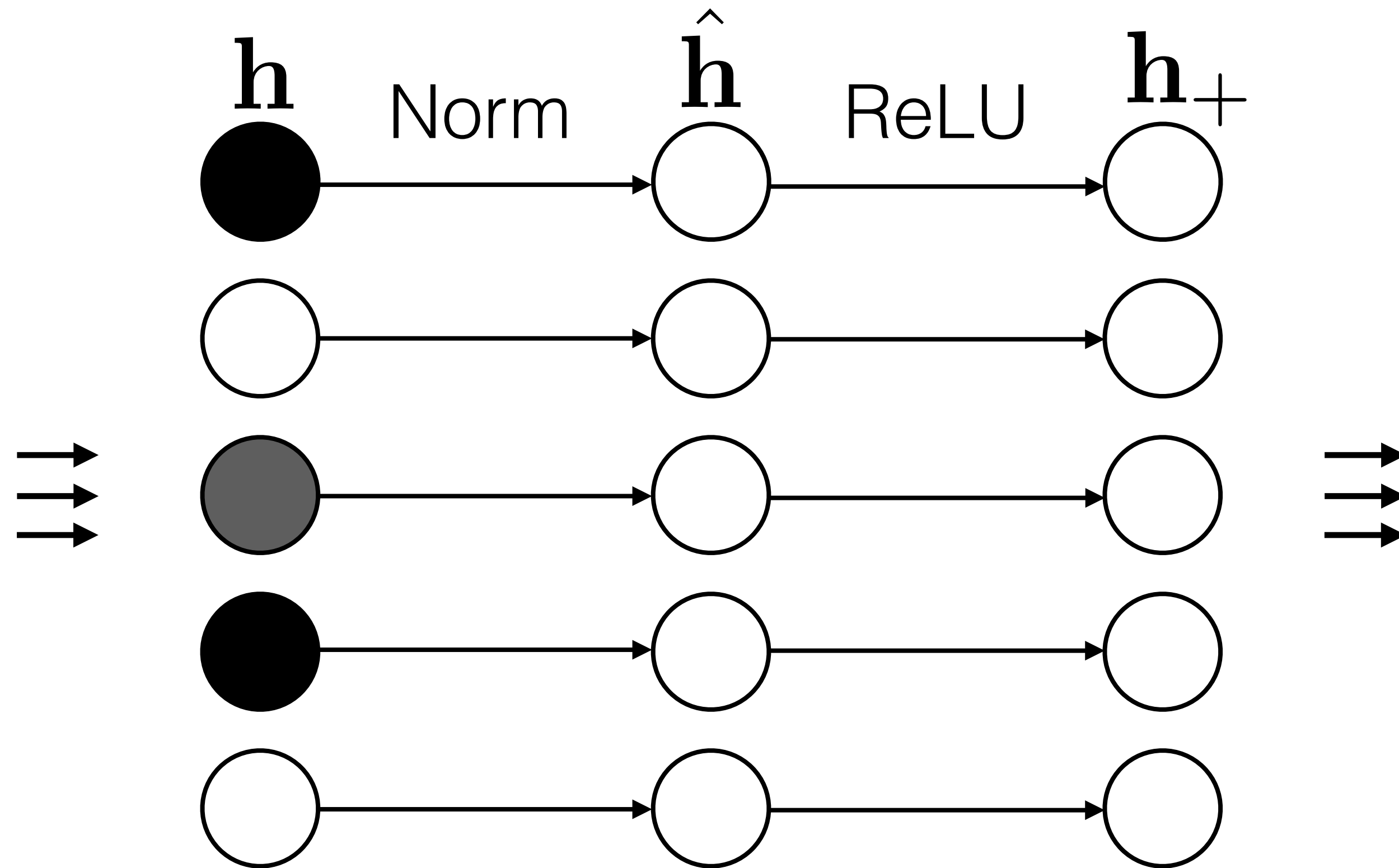
Prevents network from relying too much on spurious correlations between different hidden units.

Can be understood as averaging over an exponential **ensemble** of subnetworks. This averaging smooths the function, thereby reducing the effective capacity of the network.

Normalization layers

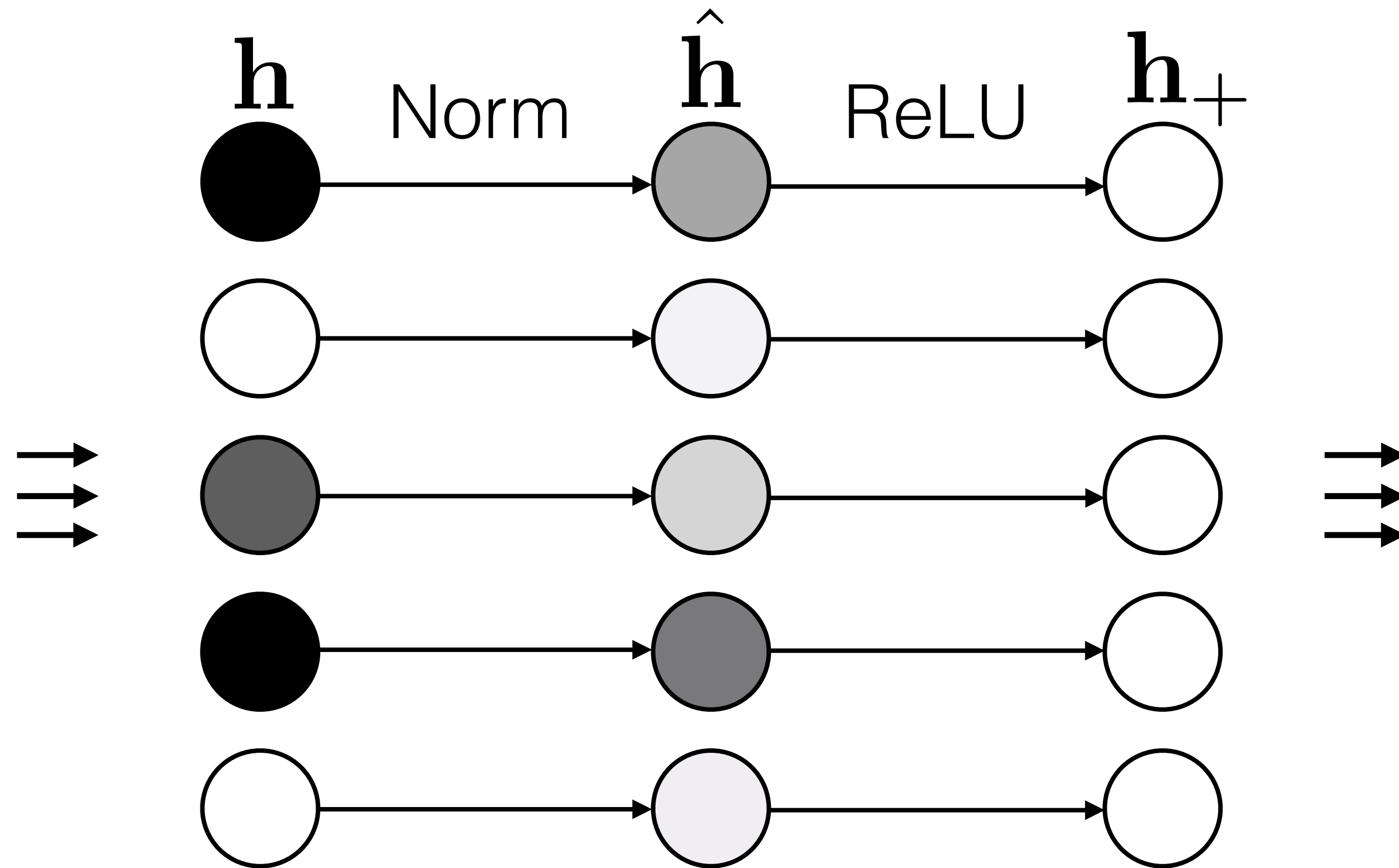


Normalization layers



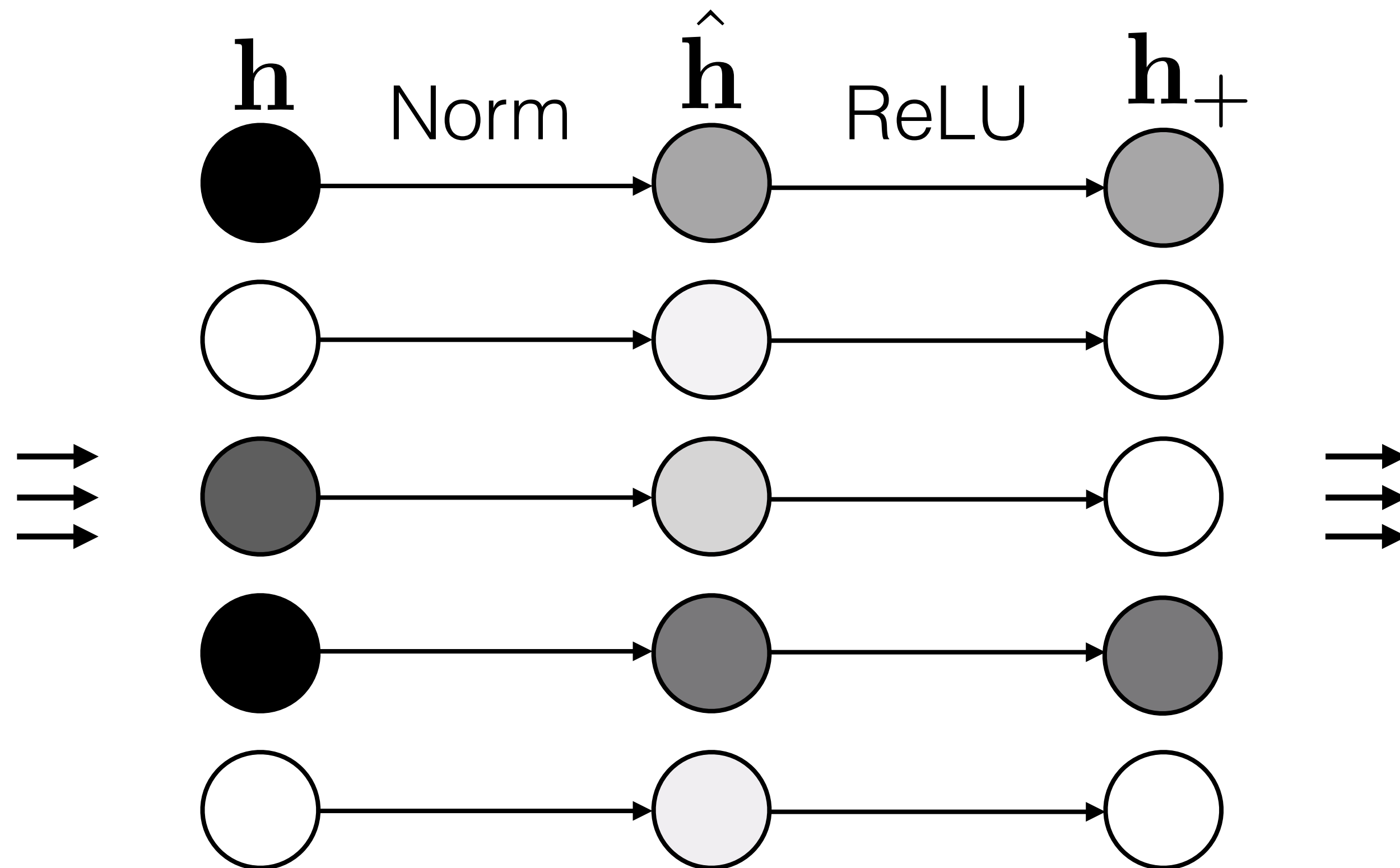
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers

Keep track of mean and variance of a unit (or a population of units) over time.

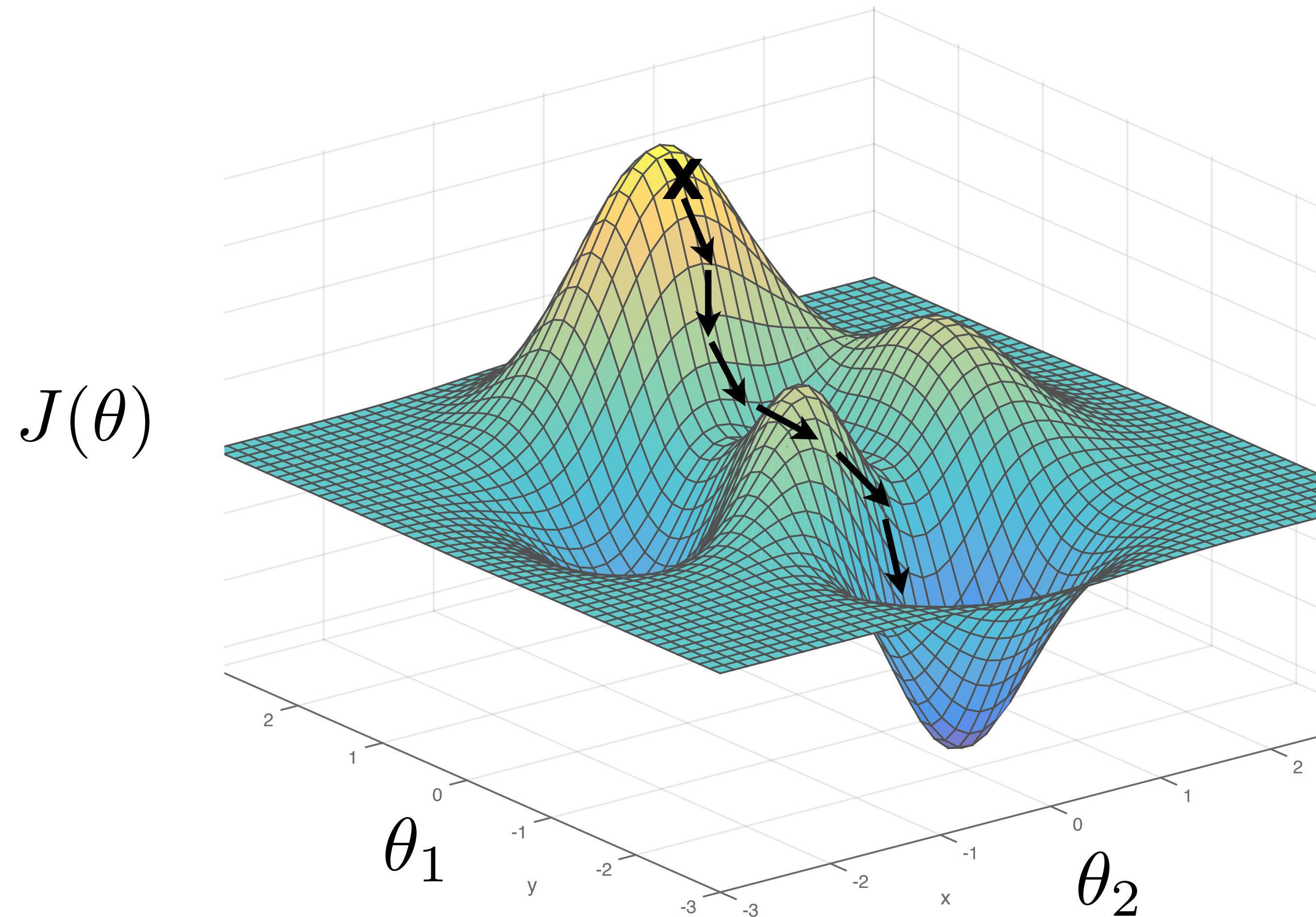
Standardize unit activations by subtracting mean and dividing by variance.

Squashes units into a **standard range**, avoiding overflow.

Also achieves **invariance** to mean and variance of the training signal.

Both these properties reduce the effective capacity of the model, i.e. regularize the model.

Gradient descent



$$\theta^* = \arg \min_{\theta} J(\theta)$$

Gradient descent

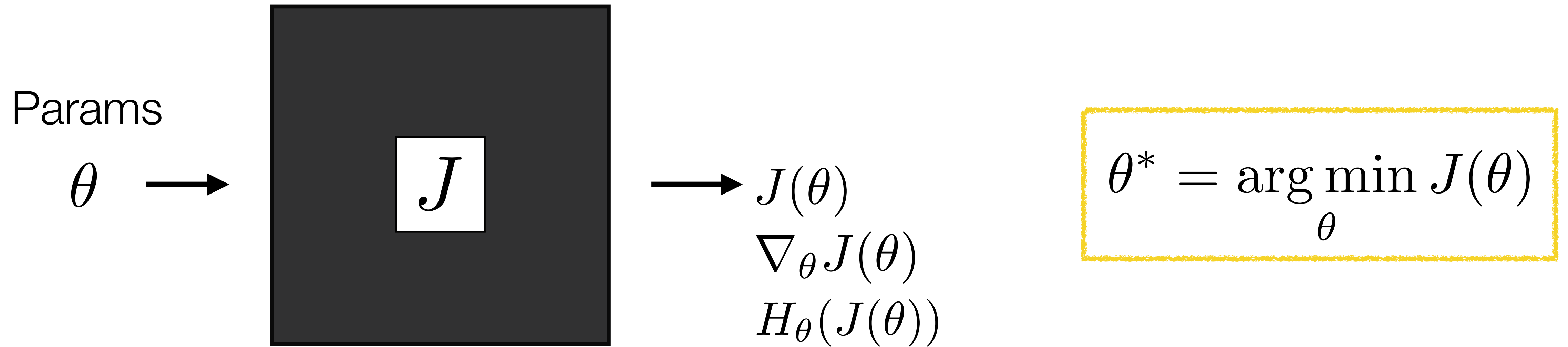
$$\theta^* = \arg \min_{\theta} \underbrace{\sum_{i=1}^N \mathcal{L}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)}_{J(\theta)}$$

One iteration of gradient descent:

$$\theta^{t+1} = \theta^t - \eta_t \left. \frac{\partial J(\theta)}{\partial \theta} \right|_{\theta=\theta^t}$$

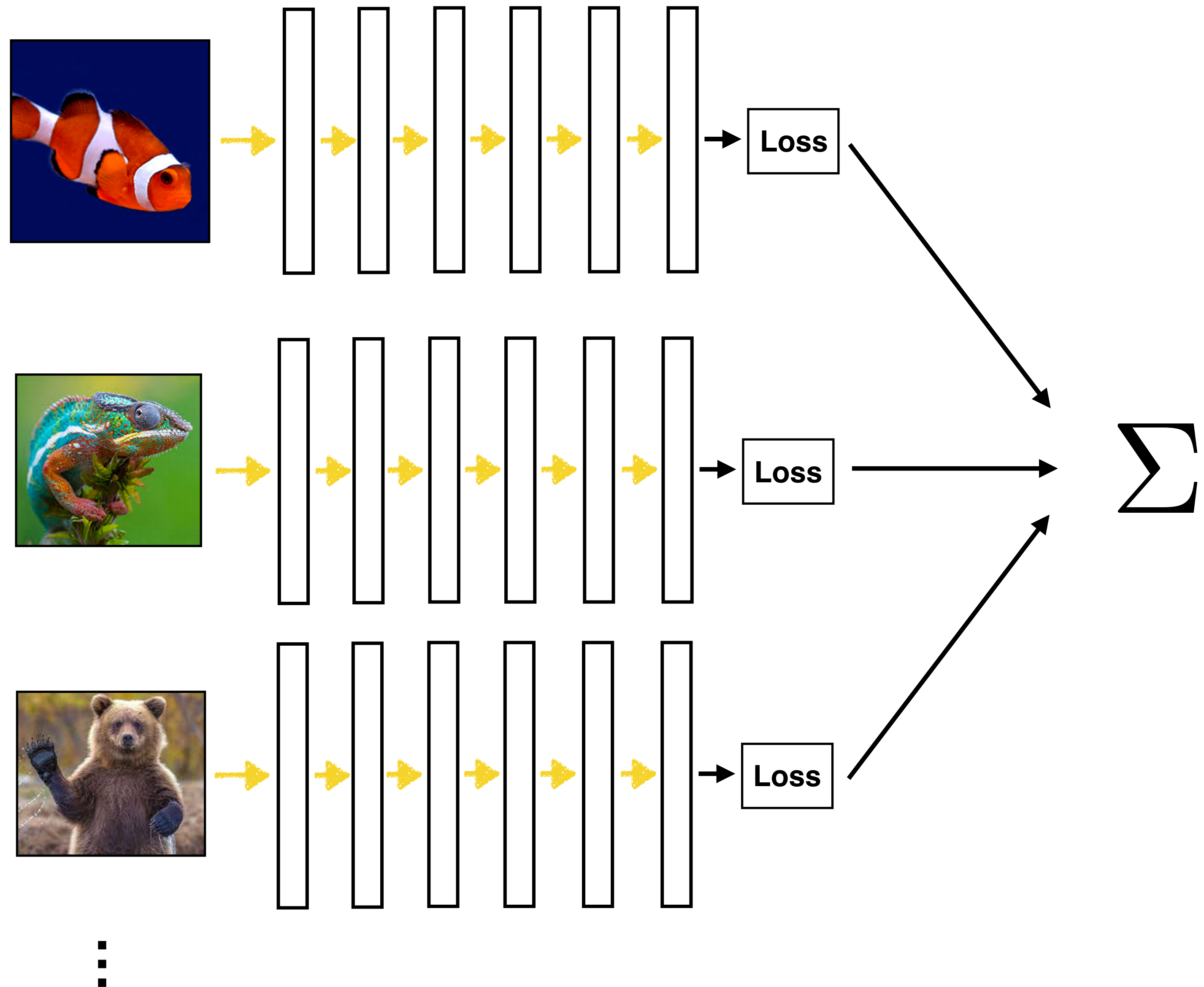
 **learning rate**

Optimization



- What's the knowledge we have about J?
 - We can evaluate $J(\theta)$
 - We can evaluate $J(\theta)$ and $\nabla_{\theta} J(\theta)$
 - We can evaluate $J(\theta)$, $\nabla_{\theta} J(\theta)$, and $H_{\theta}(J(\theta))$
- Gradient**
- Hessian**
- ← Black box optimization
- ← First order optimization
- ← Second order optimization

Batch (parallel) processing



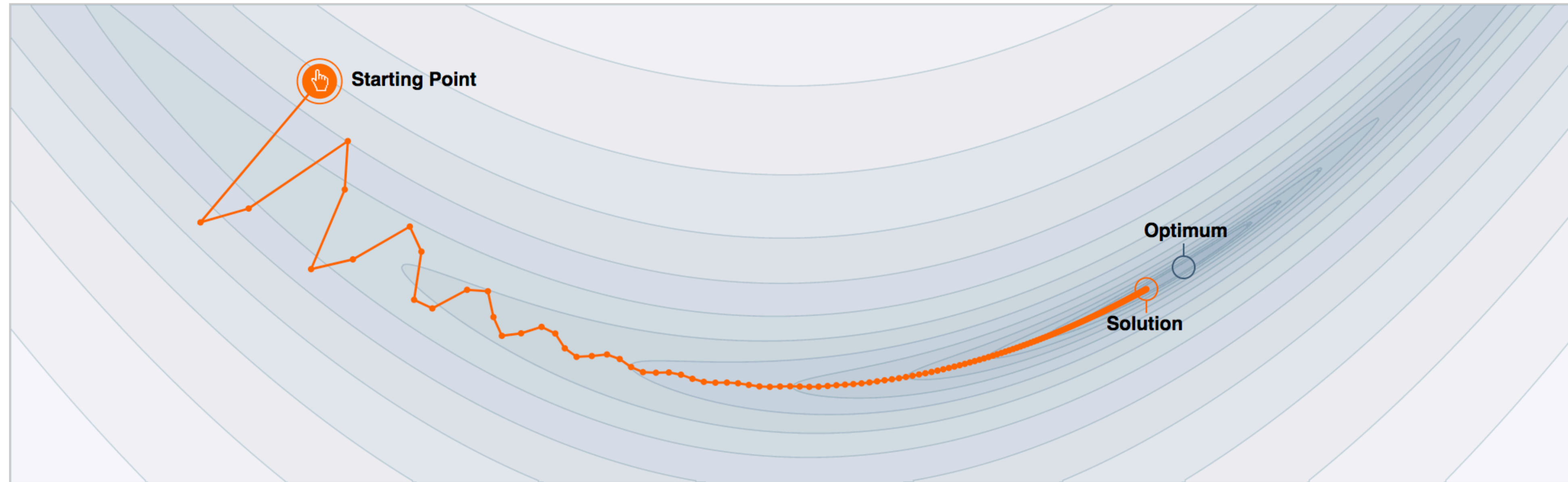
Stochastic gradient descent (SGD)

- Want to minimize overall loss function \mathbf{J} , which is sum of individual losses over each example.
- In Stochastic gradient descent, compute gradient on sub-set (batch) of data.
 - If batchsize=1 then θ is updated after each example.
 - If batchsize=N (full set) then this is standard gradient descent.
- Gradient direction is noisy, relative to average over all examples (standard gradient descent).
- Advantages
 - Faster: approximate total gradient with small sample
 - Implicit regularizer
- Disadvantages
 - High variance, unstable updates

Momentum

- Basic idea: like a ball rolling down a hill, we should build up speed so as to make faster progress when “on a roll”
- Can dampen oscillations in SGD updates
- Common in popular variants of SGD
 - Nesterov’s method
 - RMSProp
 - Adam

Why Momentum Really Works



Step-size $\alpha = 0.02$



Momentum $\beta = 0.99$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

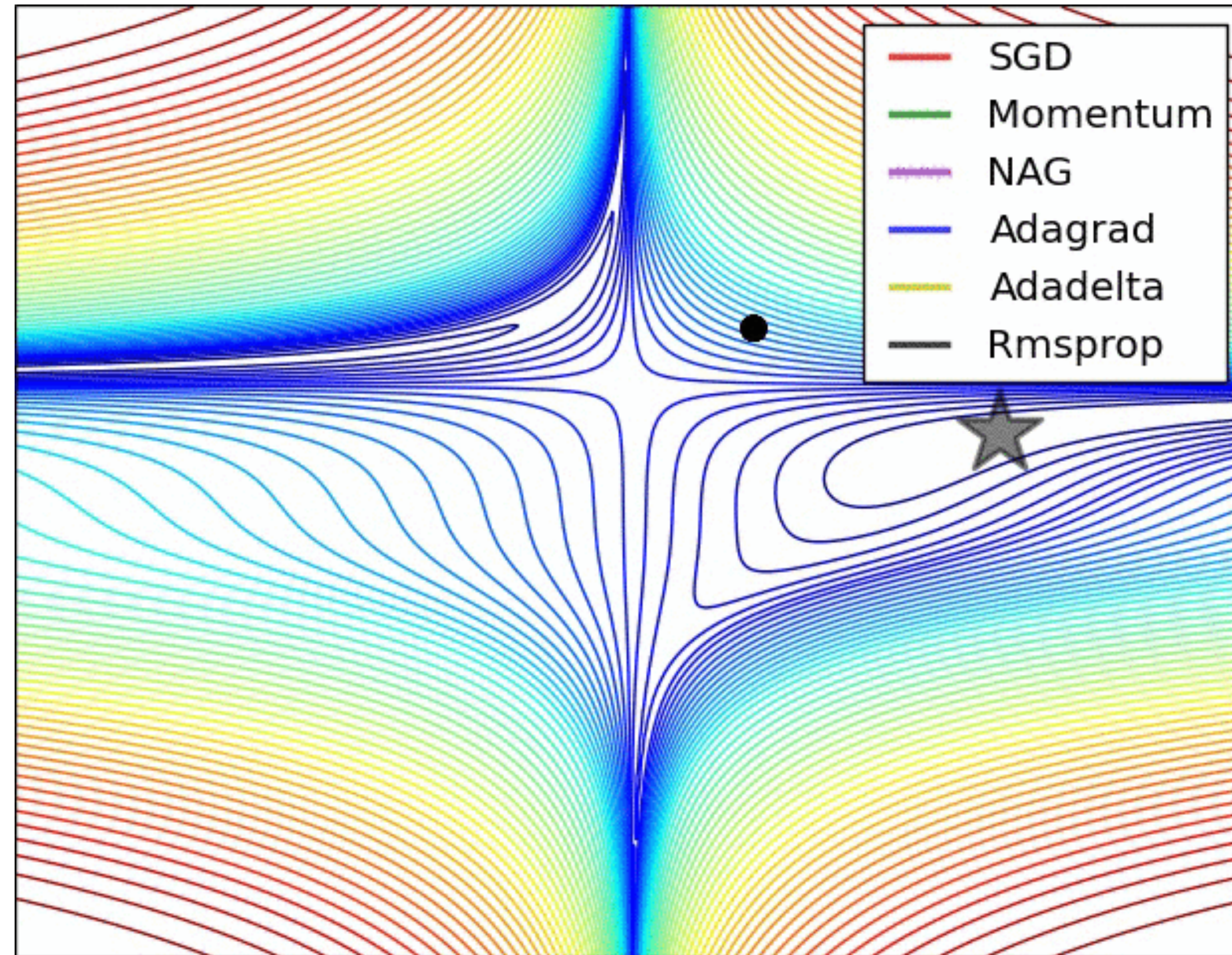
GABRIEL GOH
UC Davis

April. 4
2017

Citation:
Goh, 2017

[<https://distill.pub/2017/momentum/>]

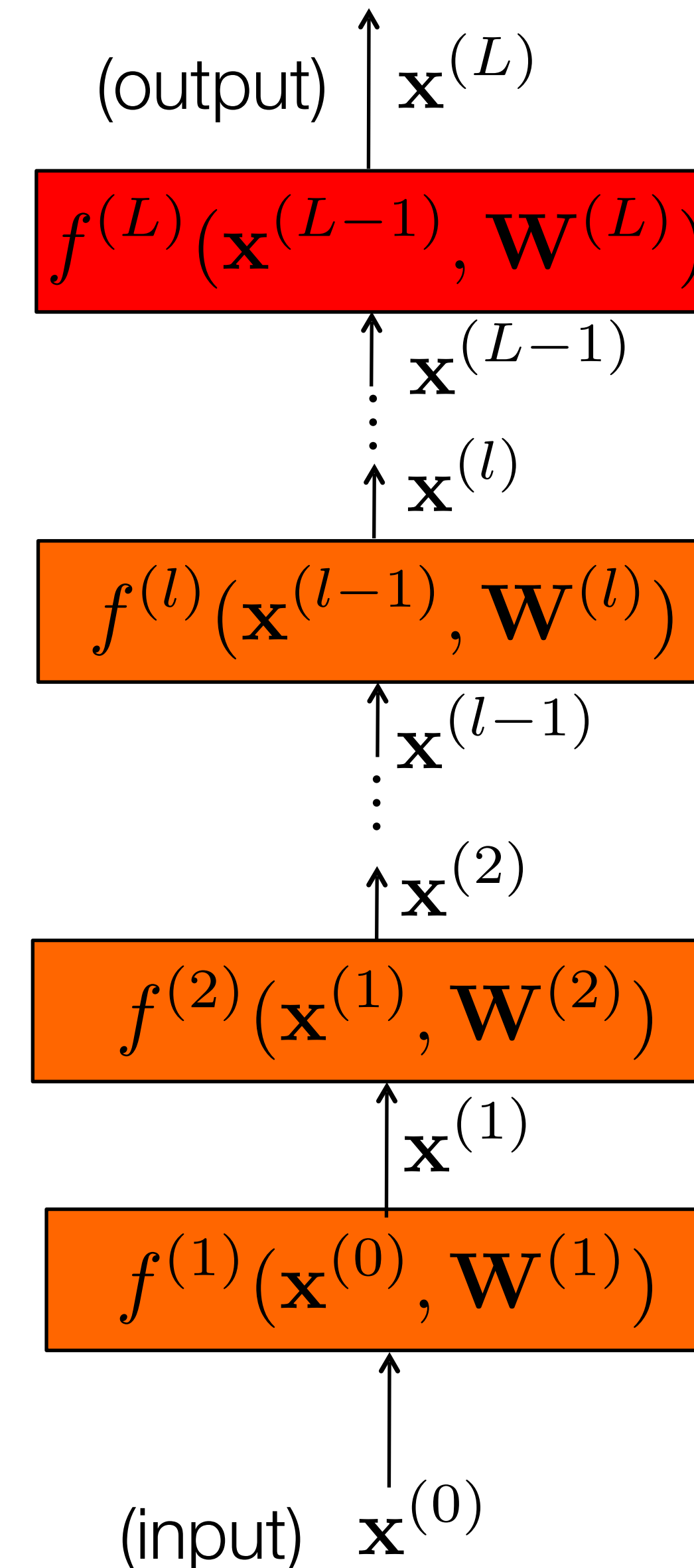
Comparison of gradient descent variants



[<http://ruder.io/optimizing-gradient-descent/>]

Forward pass

- Consider model with L layers. Layer l has vector of weights $\mathbf{W}^{(l)}$
- Forward pass:** takes input $\mathbf{x}^{(l-1)}$ and passes it through each layer $f^{(l)}$:
$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$
- Output of layer l is $\mathbf{x}^{(l)}$.
- Network output (top layer) is $\mathbf{x}^{(L)}$.

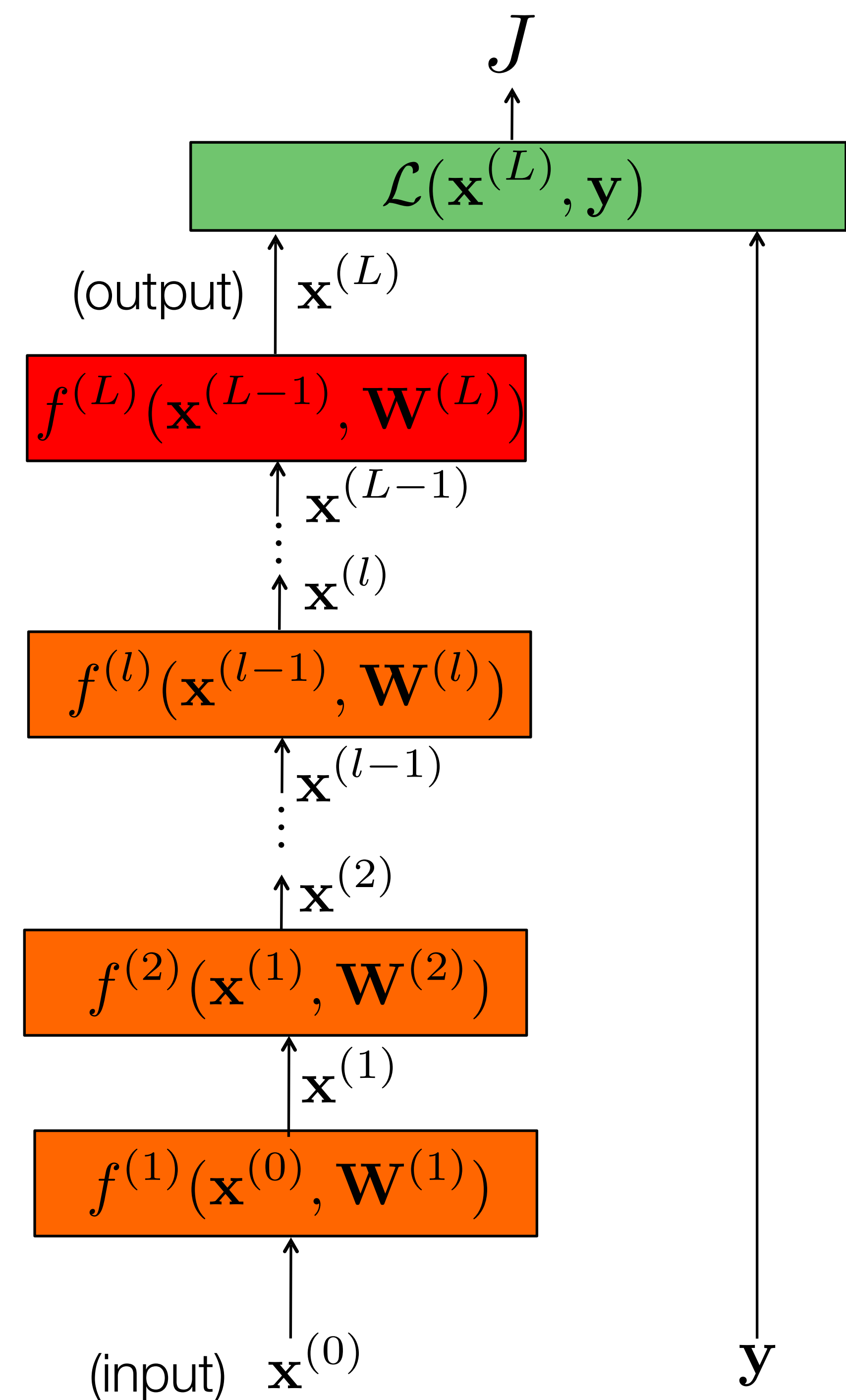


Forward pass

- Consider model with L layers. Layer l has vector of weights $\mathbf{W}^{(l)}$
- Forward pass:** takes input $\mathbf{x}^{(l-1)}$ and passes it through each layer $f^{(l)}$:

$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

- Output of layer l is $\mathbf{x}^{(l)}$.
- Network output (top layer) is $\mathbf{x}^{(L)}$.
- Loss function** \mathcal{L} compares $\mathbf{x}^{(L)}$ to \mathbf{y} .
- Overall energy is the sum of the cost over all training examples: $J = \sum_{i=1}^N \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i)$



Gradient descent

- We need to compute gradients of the cost with respect to model parameters $\mathbf{W}^{(l)}$.
- By design, each layer is differentiable with respect to its parameters and input.

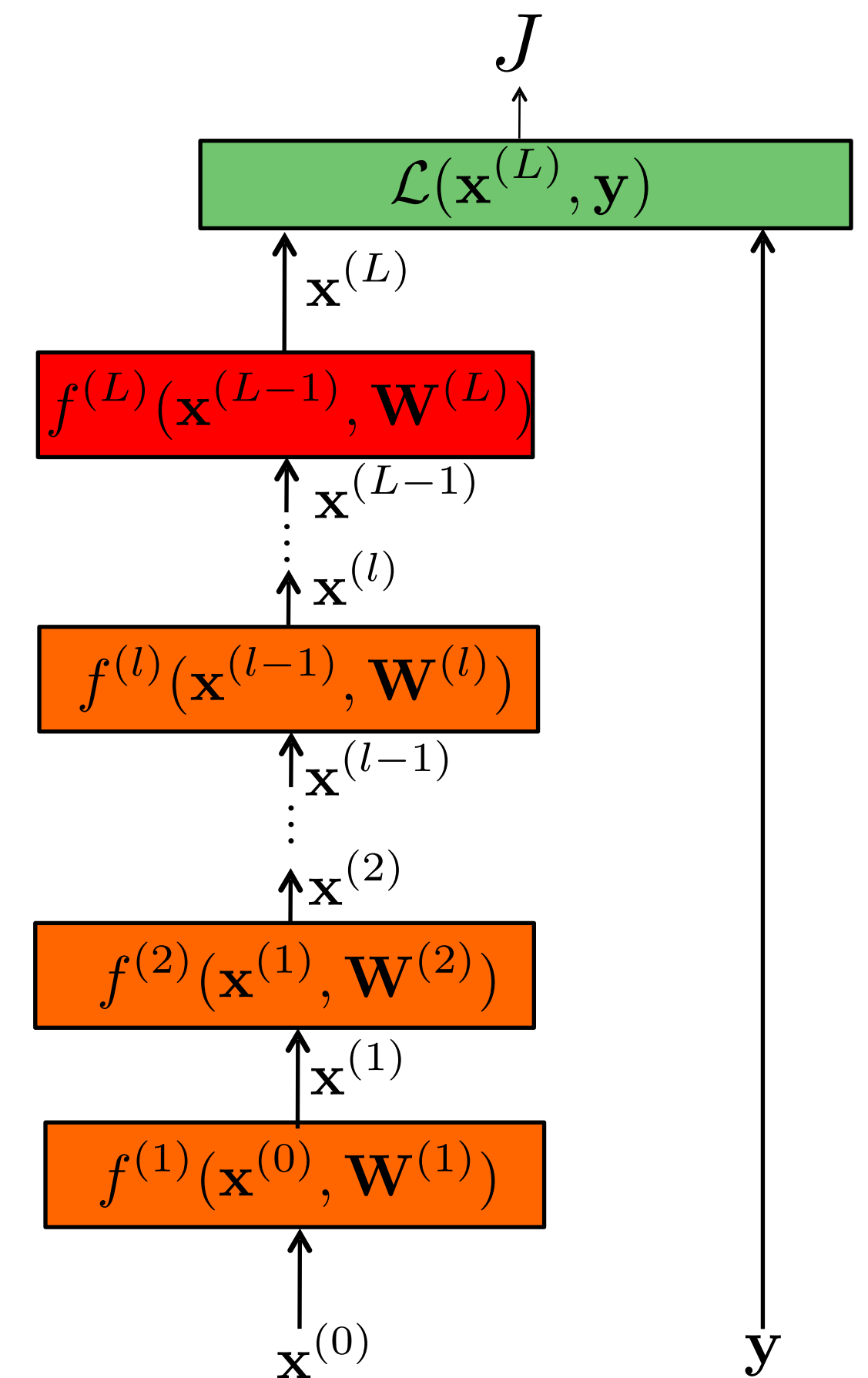
Computing gradients

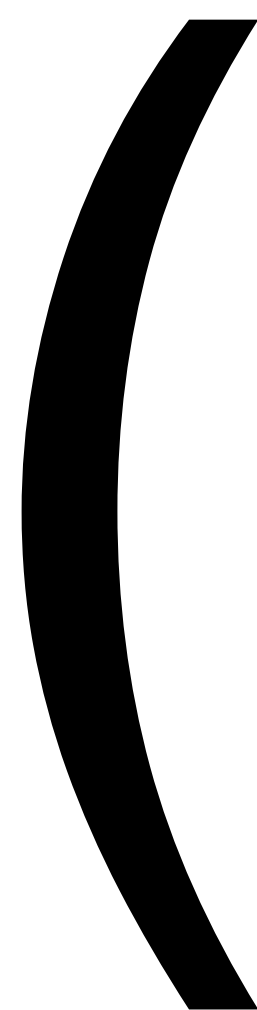
To compute the gradients, we could start by writing the full energy J as a function of the network parameters.

$$J(\mathbf{W}) = \sum_{i=1} \mathcal{L}(f^{(L)}(\dots f^{(2)}(f^{(1)}(\mathbf{x}_i^{(0)}, \mathbf{W}^{(1)}), \mathbf{W}^{(2)}), \dots, \mathbf{W}^{(L)}), \mathbf{y}_i)$$

And then compute the partial derivatives... instead, we can use the chain rule to derive a compact algorithm:

backpropagation





Matrix calculus

- \mathbf{x} column vector of size $[n \times 1]$:
$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

- We now define a function on vector \mathbf{x} : $\mathbf{y} = f(\mathbf{x})$
- If y is a scalar, then

$$\frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right)$$

The derivative of y is a row vector of size $[1 \times n]$

- If \mathbf{y} is a vector $[m \times 1]$, then (*Jacobian formulation*):

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{pmatrix}$$

The derivative of \mathbf{y} is a matrix of size $[m \times n]$

(m rows and n columns)

Matrix calculus

- If y is a scalar and \mathbf{X} is a matrix of size $[n \times m]$, then

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \cdots & \frac{\partial y}{\partial x_{n1}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y}{\partial x_{1m}} & \frac{\partial y}{\partial x_{2m}} & \cdots & \frac{\partial y}{\partial x_{nm}} \end{pmatrix}$$

The output is a matrix of size $[m \times n]$

Matrix calculus

- Chain rule:

For the function: $\mathbf{z} = h(\mathbf{x}) = f(g(\mathbf{x}))$

Its derivative is: $h'(\mathbf{x}) = f'(g(\mathbf{x}))g'(\mathbf{x})$

and writing $\mathbf{z} = f(\mathbf{u})$, and $\mathbf{u} = g(\mathbf{x})$:

$$\left. \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{a}} = \left. \frac{\partial \mathbf{z}}{\partial \mathbf{u}} \right|_{\mathbf{u}=g(\mathbf{a})} \cdot \left. \frac{\partial \mathbf{u}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{a}}$$

\uparrow
 $[m \times n]$

\uparrow
 $[m \times p]$

\uparrow
 $[p \times n]$

with $p = \text{length of vector } \mathbf{u} = |\mathbf{u}|$, $m = |\mathbf{z}|$, and $n = |\mathbf{x}|$

Example, if $|\mathbf{z}| = 1$, $|\mathbf{u}| = 2$, $|\mathbf{x}| = 4$

$$h'(\mathbf{x}) = \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} = \begin{array}{|c|c|} \hline \text{blue} & \text{blue} \\ \hline \end{array} \begin{array}{|c|c|c|c|} \hline \text{red} & \text{red} & \text{red} & \text{red} \\ \hline \text{red} & \text{red} & \text{red} & \text{red} \\ \hline \end{array}$$

Matrix calculus

- Chain rule:

For the function: $h(\mathbf{x}) = f^{(n)}(f^{(n-1)}(\dots f^{(1)}(\mathbf{x})))$

With $\mathbf{u}^{(1)} = f^{(1)}(\mathbf{x})$

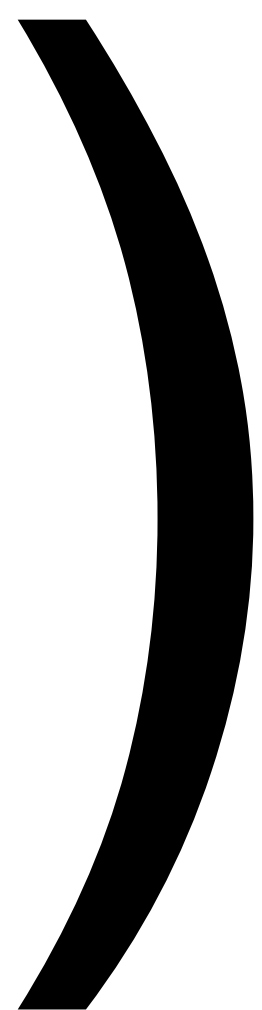
$$\mathbf{u}^{(i)} = f^{(i)}(\mathbf{u}^{(i-1)})$$

$$\mathbf{z} = \mathbf{u}^{(n)} = f^{(n)}(\mathbf{u}^{(n-1)})$$

The derivative becomes a product of matrices:

$$\left. \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{a}} = \left. \frac{\partial \mathbf{z}}{\partial \mathbf{u}^{(n-1)}} \right|_{\mathbf{u}^{(n-1)}=f^{(n-1)}(\mathbf{u}^{(n-2)})} \cdot \left. \frac{\partial \mathbf{u}^{(n-1)}}{\partial \mathbf{u}^{(n-2)}} \right|_{\mathbf{u}^{(n-2)}=f^{(n-2)}(\mathbf{u}^{(n-3)})} \cdots \left. \frac{\partial \mathbf{u}^{(2)}}{\partial \mathbf{u}^{(1)}} \right|_{\mathbf{u}^{(1)}=f^{(1)}(\mathbf{a})} \cdot \left. \frac{\partial \mathbf{u}^{(1)}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{a}}$$

(exercise: check that all the matrix dimensions work out fine)



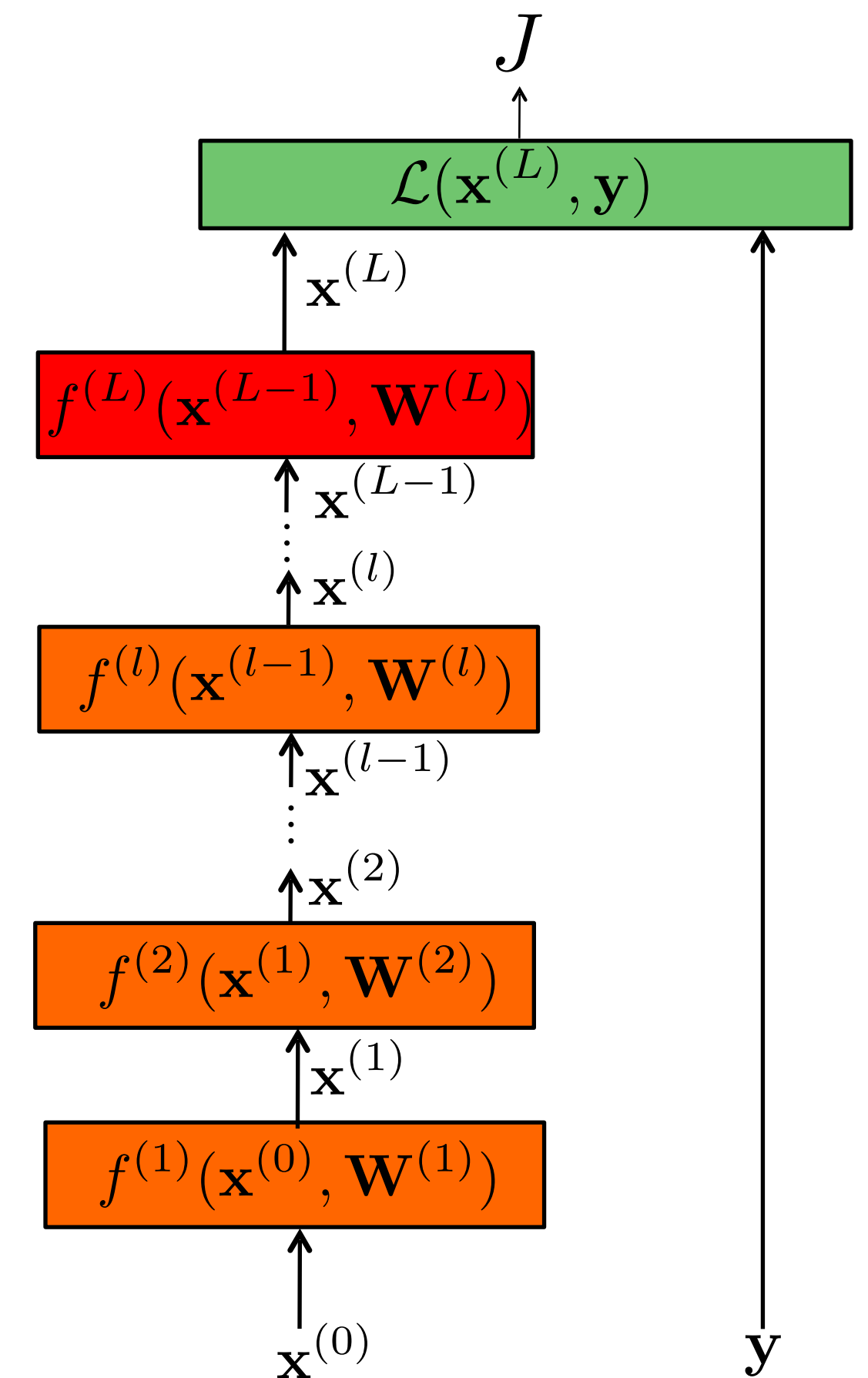
Computing gradients

To compute the gradients, we could start by writing the full energy J as a function of the network parameters.

$$J(\mathbf{W}) = \sum_{i=1} \mathcal{L}(f^{(L)}(\dots f^{(2)}(f^{(1)}(\mathbf{x}_i^{(0)}, \mathbf{W}^{(1)}), \mathbf{W}^{(2)}), \dots \mathbf{W}^{(L)}), \mathbf{y}_i)$$

And then compute the partial derivatives... instead, we can use the chain rule to derive a compact algorithm:

backpropagation



Computing gradients

The energy J is the sum of the losses associated to each training example $\{\mathbf{x}_i^{(0)}, \mathbf{y}_i\}$

$$J(\mathbf{W}) = \sum_{i=1}^N \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i; \mathbf{W})$$

Its gradient with respect to each of the network's parameters w is:

$$\frac{\partial J(\mathbf{W})}{\partial w} = \sum_{i=1}^N \frac{\partial \mathcal{L}(\mathbf{x}_i^{(L)}, \mathbf{y}_i; \mathbf{W})}{\partial w}$$

is how much J varies when the parameter w is varied.

Computing gradients

We could write the loss function to get the gradients as:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}; \mathbf{W}) = \mathcal{L}(f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)}), \mathbf{y})$$


If we compute the gradient with respect to the parameters of the last layer (output layer) $\mathbf{W}^{(L)}$, using the **chain rule**:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})}{\partial \mathbf{W}^{(L)}}$$

(How much the cost changes when we change $\mathbf{W}^{(L)}$ is the product between how much the loss changes when we change the output of the last layer and how much the output changes when we change the layer parameters.)

Computing gradients: loss layer

If we compute the gradient with respect to the parameters of the last layer (output layer) $\mathbf{W}^{(L)}$, using the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{W}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial f^{(L)}(\mathbf{x}^{(L-1)}, \mathbf{W}^{(L)})}{\partial \mathbf{W}^{(L)}}$$


For example, for an Euclidean loss:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}) = \frac{1}{2} \left\| \mathbf{x}^{(L)} - \mathbf{y} \right\|_2^2$$

Will depend on the layer structure and non-linearity.

The gradient is:

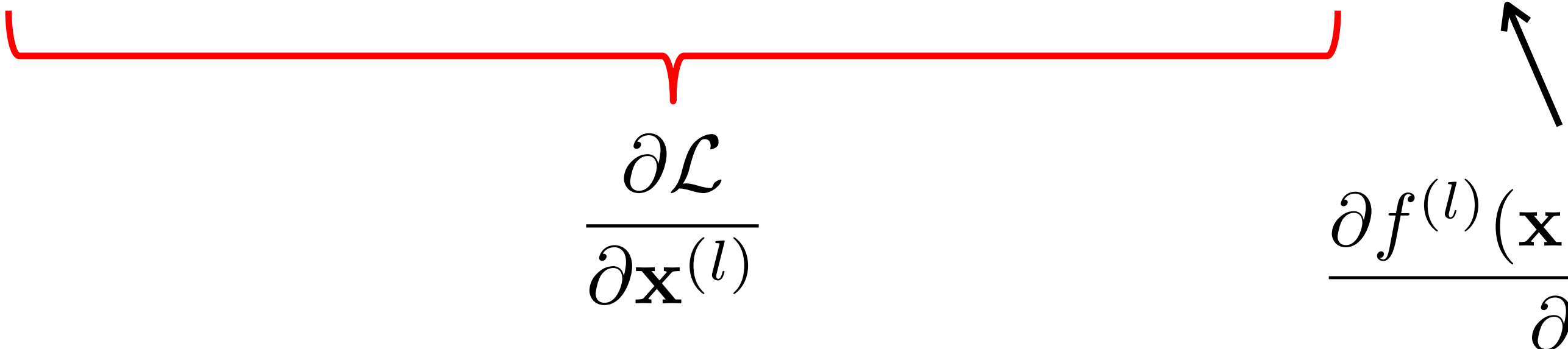
$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} = \mathbf{x}^{(L)} - \mathbf{y}$$

Computing gradients: layer l

We could write the full loss function to get the gradients:

$$\mathcal{L}(\mathbf{x}^{(L)}, \mathbf{y}; \mathbf{W}) = \mathcal{L}(f^{(L)}(\dots f^{(2)}(f^{(1)}(\mathbf{x}^{(0)}, \mathbf{W}^{(1)}), \mathbf{W}^{(2)}), \dots, \mathbf{W}^{(L)}), \mathbf{y})$$

If we compute the gradient with respect to w_i , using the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \underbrace{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}}}_{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}$$


And this can be
computed iteratively!

This is easy.

Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(L)}} = \underbrace{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \cdot \frac{\partial \mathbf{x}^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(l+1)}}{\partial \mathbf{x}^{(l)}}}_{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{W}^{(l)}}$$

\nearrow
 $\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$

If we have the value of $\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}}$ we can compute the gradient at the layer below as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial \mathbf{x}^{(l)}}{\partial \mathbf{x}^{(l-1)}}$$

\nearrow
 Gradient
layer l-1

\nearrow
 Gradient
layer l

\nwarrow
 $\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$

Backpropagation — Goal: to update parameters of layer l

- Layer l has two inputs (during training)

$$\begin{array}{c} \mathbf{x}^{(l-1)} \xrightarrow{\text{green}} \text{orange box} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \xrightarrow{\text{red}} \text{orange box} \end{array}$$

- We compute the outputs

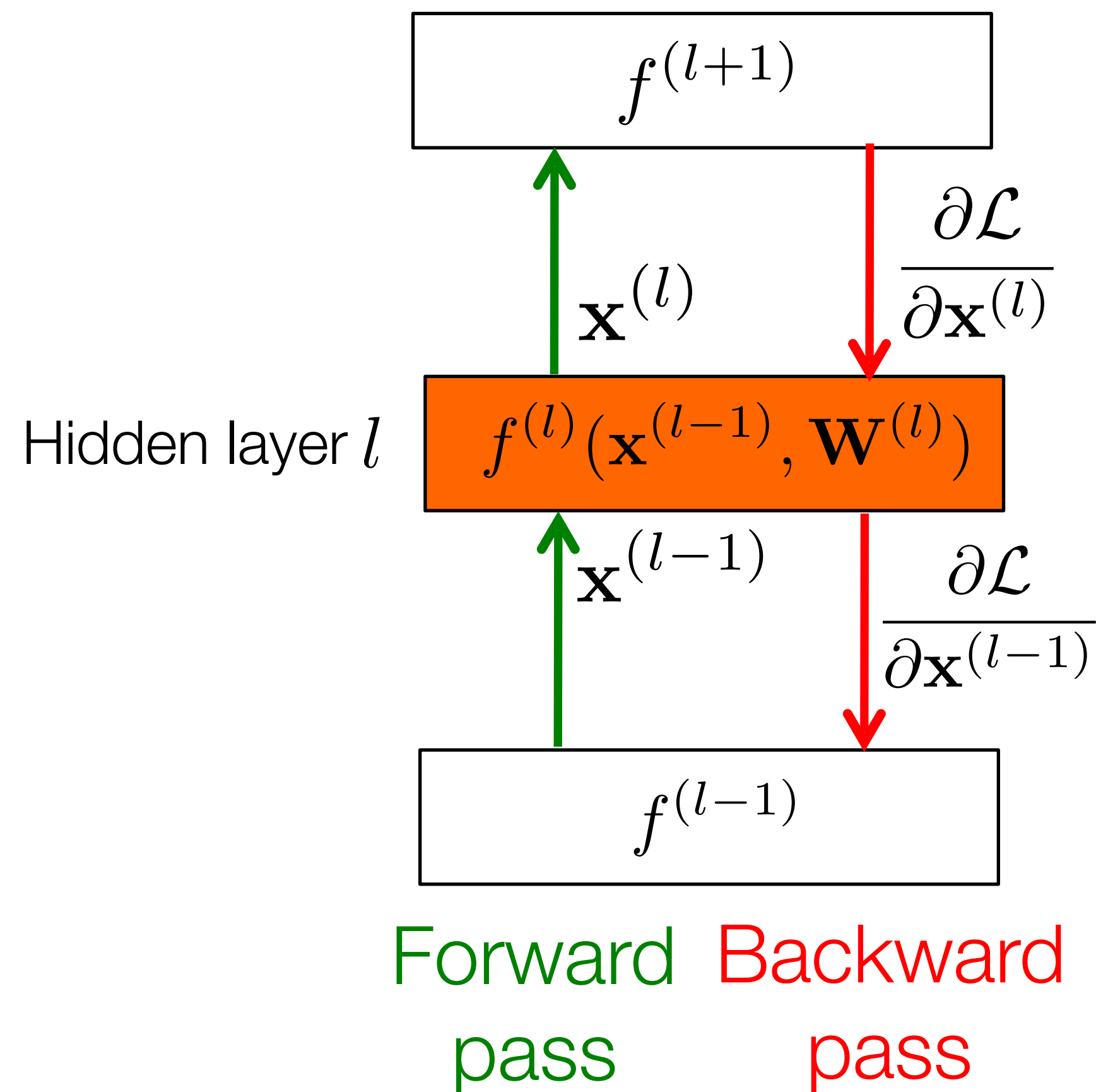
$$\begin{array}{c} \text{orange box} \xrightarrow{\text{green}} \mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)}) \\ \text{orange box} \xrightarrow{\text{red}} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}} \end{array}$$

- To compute the output, we need:

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- To compute the weight update, we need:

$$\frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$



Backpropagation — Goal: to update parameters of layer l

- Layer l has two inputs (during training)

$$\mathbf{x}^{(l-1)} \xrightarrow{\text{green}} \text{orange box}$$
$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \xrightarrow{\text{red}} \text{orange box}$$

- We compute the outputs

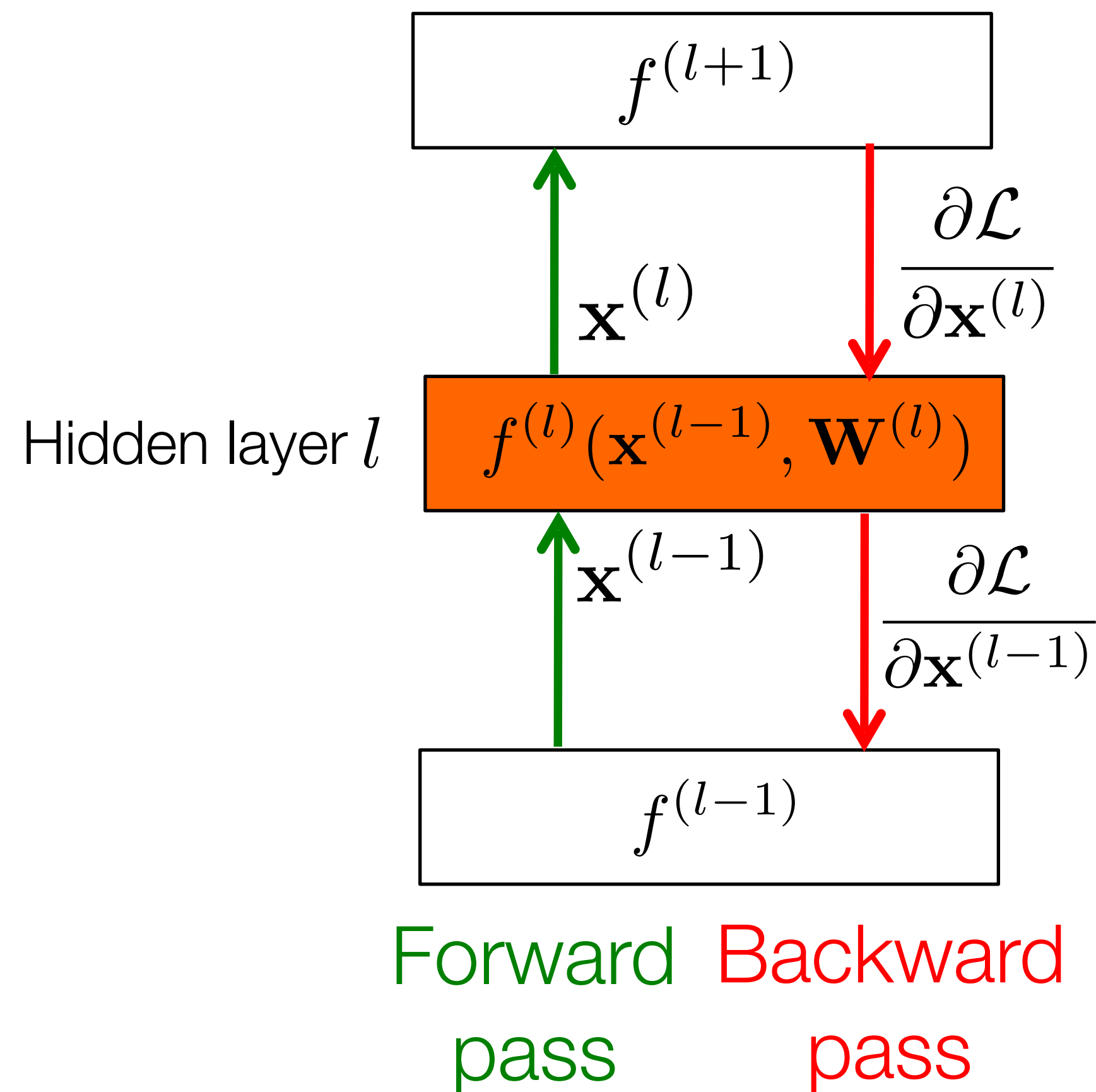
$$\text{orange box} \xrightarrow{\text{green}} \mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$
$$\text{orange box} \xrightarrow{\text{red}} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- The weight update equation is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} + \eta \left(\frac{\partial J}{\partial \mathbf{W}^{(l)}} \right)^T$$

(sum over all training examples to get J)



Backpropagation Summary

- Forward pass: for each training example, compute the outputs for all layers:

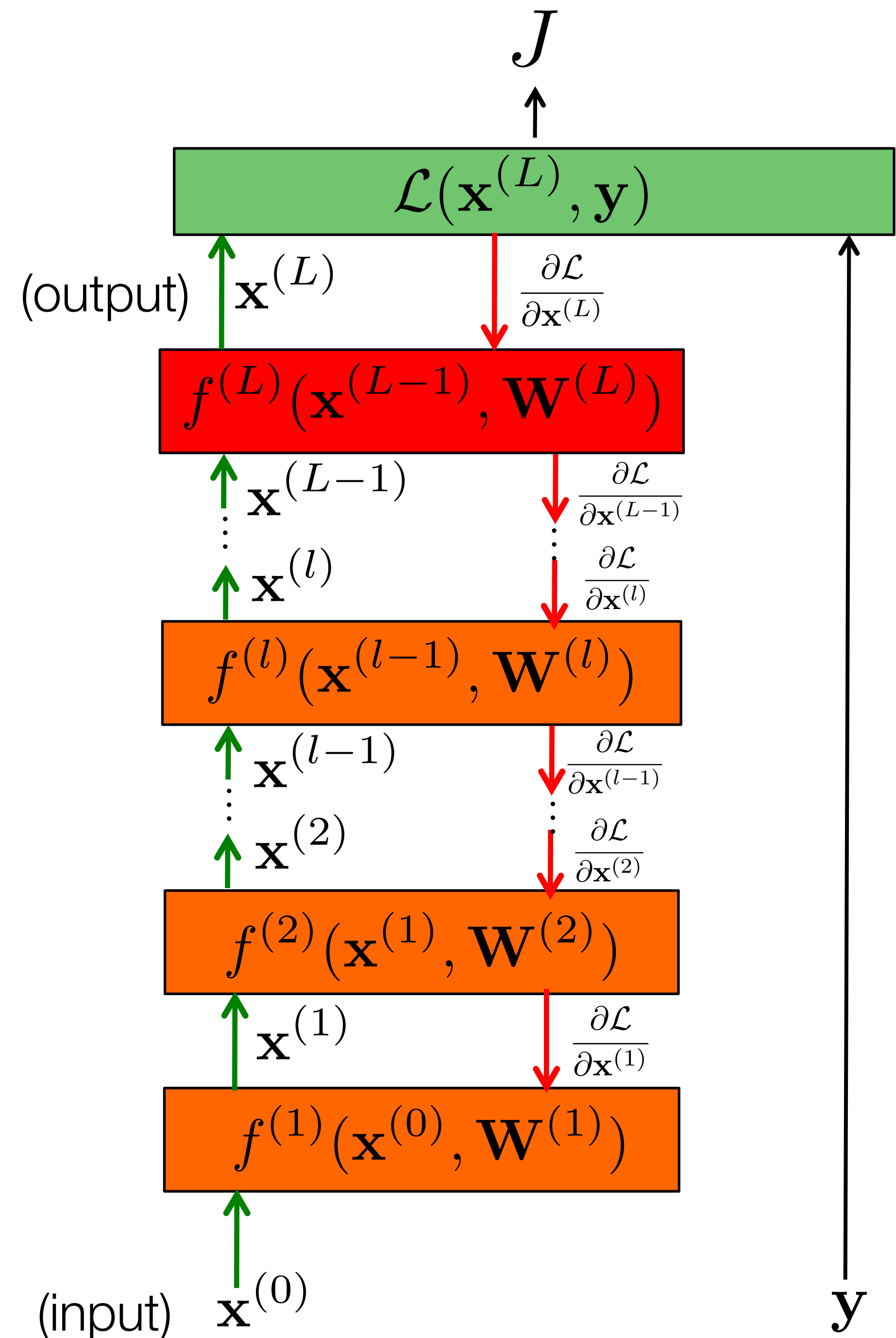
$$\mathbf{x}^{(l)} = f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})$$

- Backwards pass: compute loss derivatives iteratively from top to bottom:

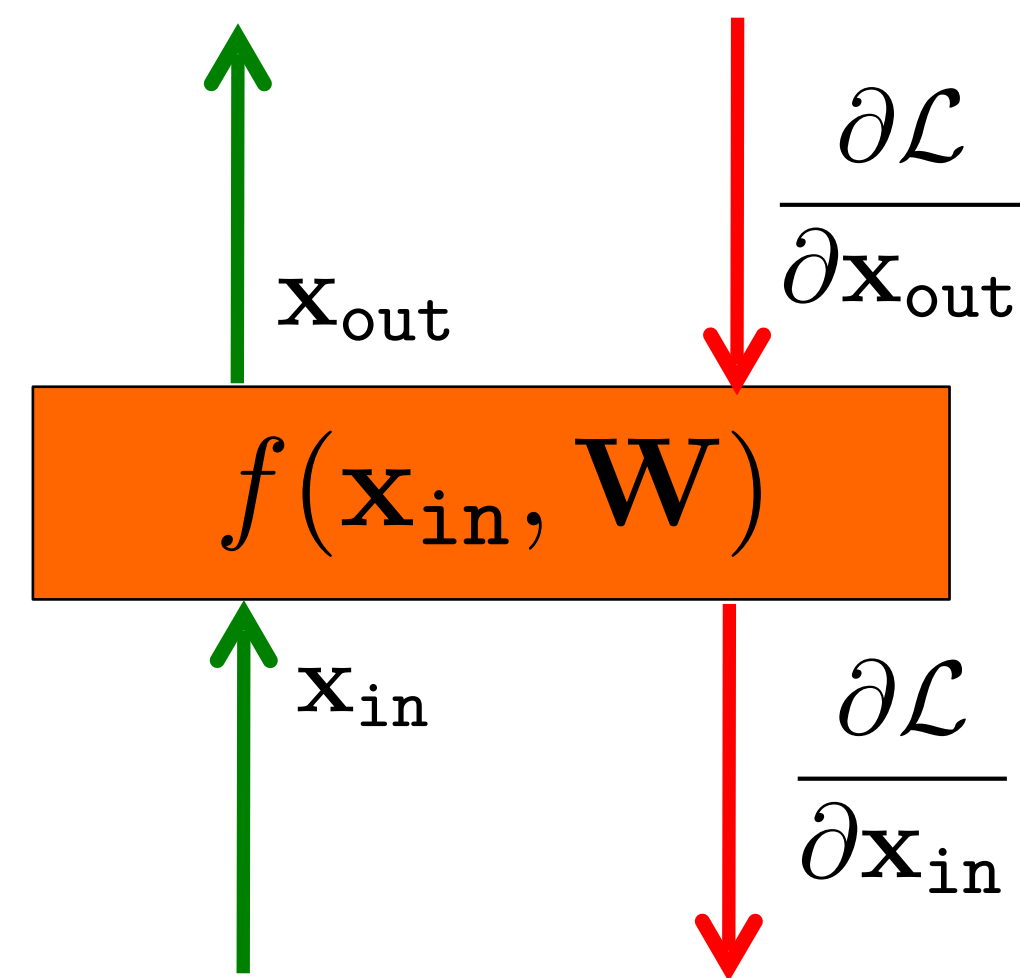
$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{x}^{(l-1)}}$$

- Compute gradients w.r.t. weights, and update weights:

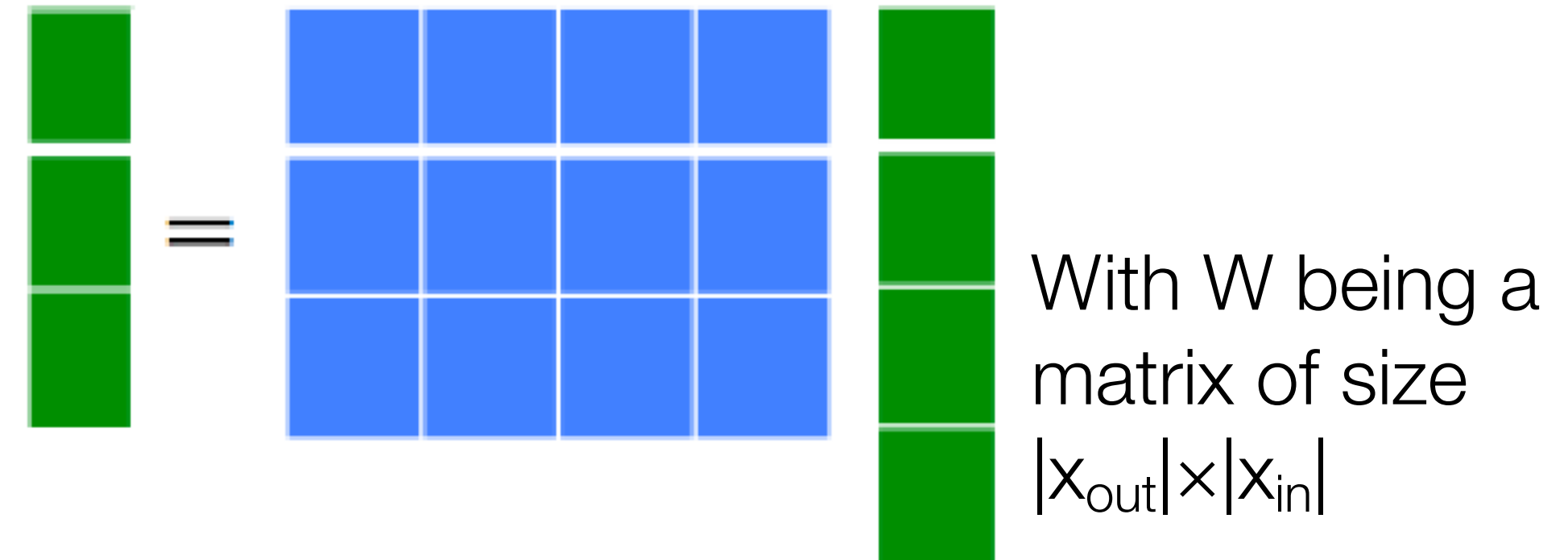
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(l)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(l)}} \cdot \frac{\partial f^{(l)}(\mathbf{x}^{(l-1)}, \mathbf{W}^{(l)})}{\partial \mathbf{W}^{(l)}}$$



Linear Module



- Forward propagation: $\mathbf{x}_{\text{out}} = f(\mathbf{x}_{\text{in}}, \mathbf{W}) = \mathbf{W}\mathbf{x}_{\text{in}}$



- Backprop to input:

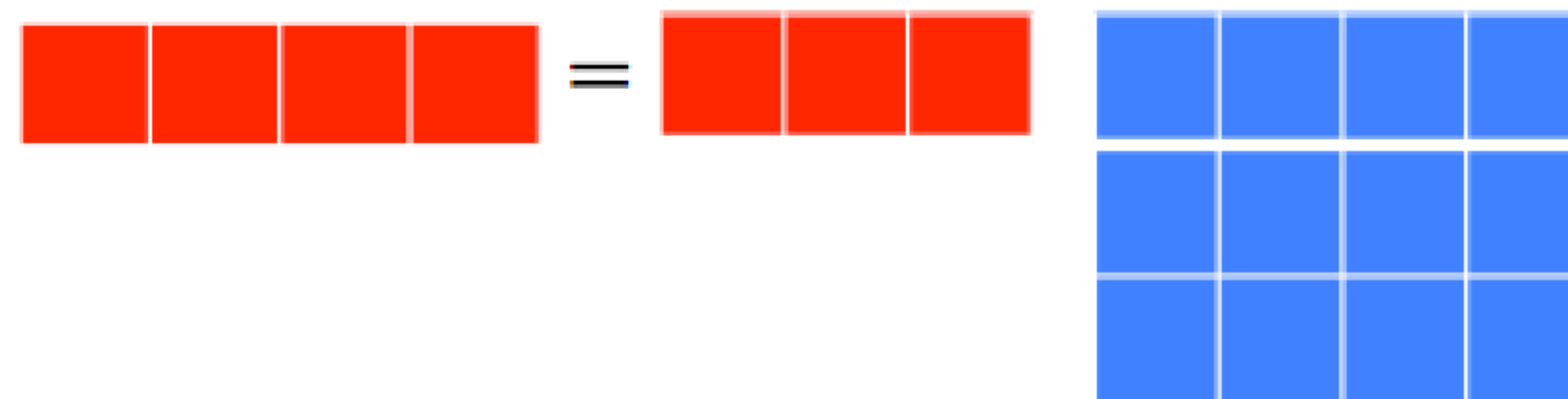
$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{in}}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}} \cdot \frac{\partial f(\mathbf{x}_{\text{in}}, \mathbf{W})}{\partial \mathbf{x}_{\text{in}}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}} \cdot \frac{\partial \mathbf{x}_{\text{out}}}{\partial \mathbf{x}_{\text{in}}}$$

If we look at the j component of output \mathbf{x}_{out} , with respect to the i component of the input, \mathbf{x}_{in} :

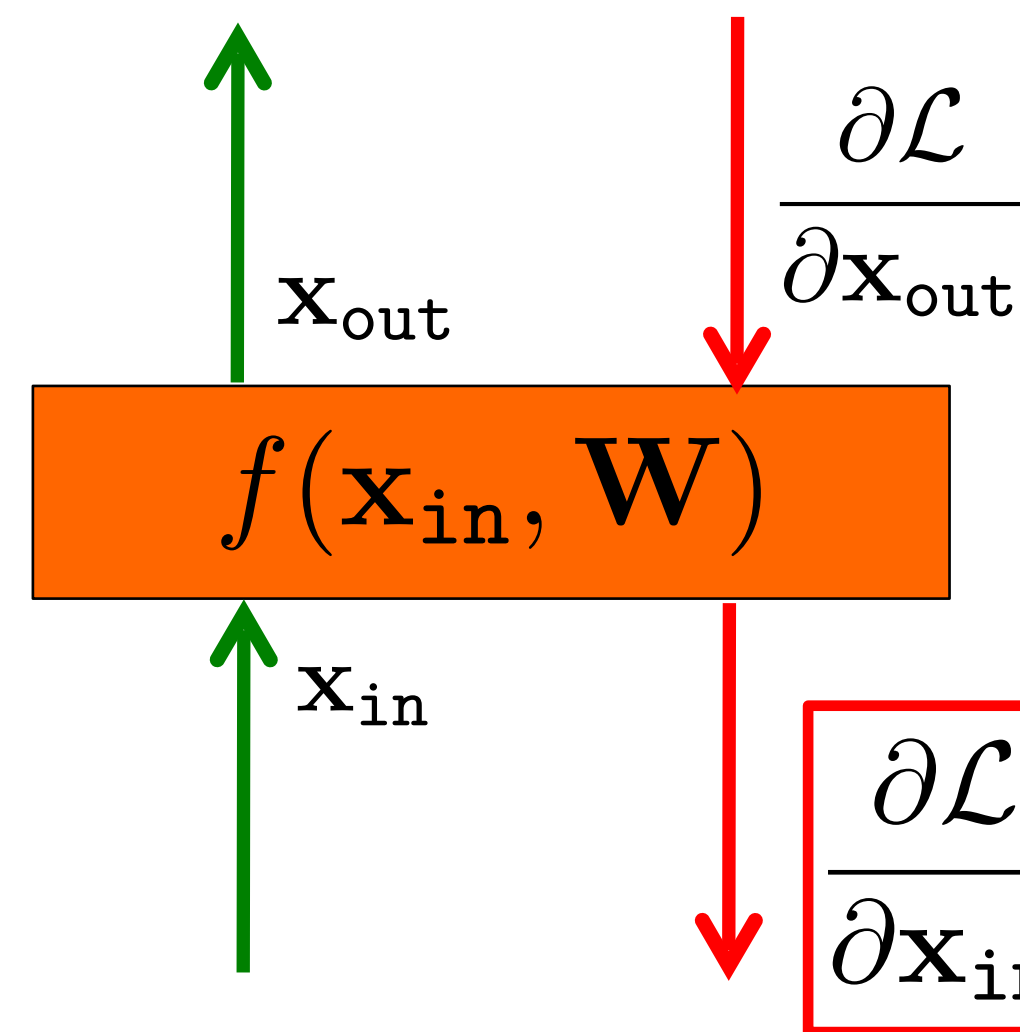
$$\frac{\partial \mathbf{x}_{\text{out}_i}}{\partial \mathbf{x}_{\text{in}_j}} = \mathbf{W}_{ij} \longrightarrow \frac{\partial f(\mathbf{x}_{\text{in}}, \mathbf{W})}{\partial \mathbf{x}_{\text{in}}} = \mathbf{W}$$

Therefore:

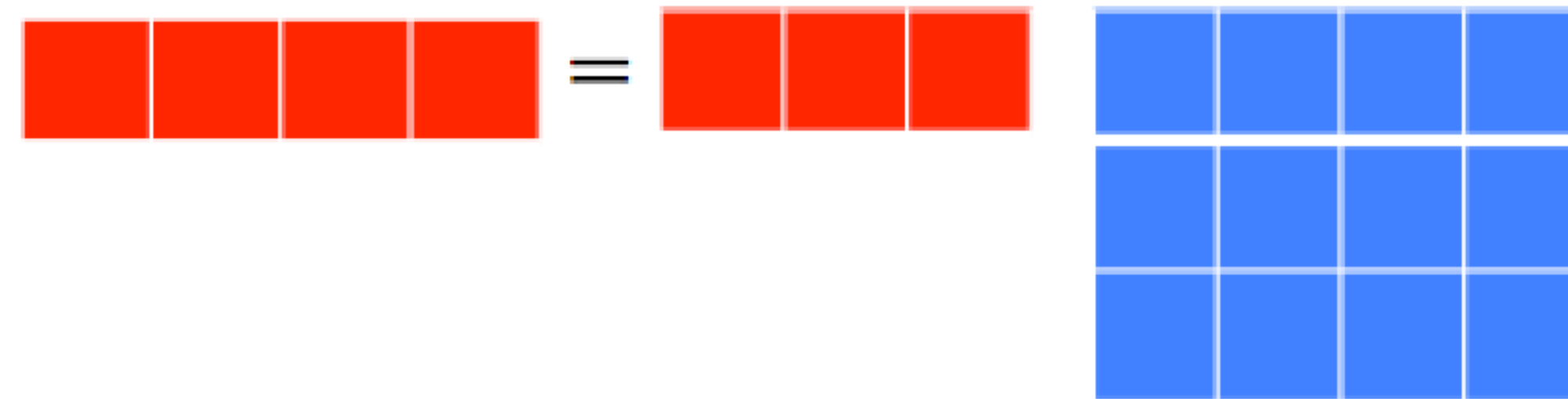
$$\boxed{\frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{in}}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}} \cdot \mathbf{W}}$$



Linear Module

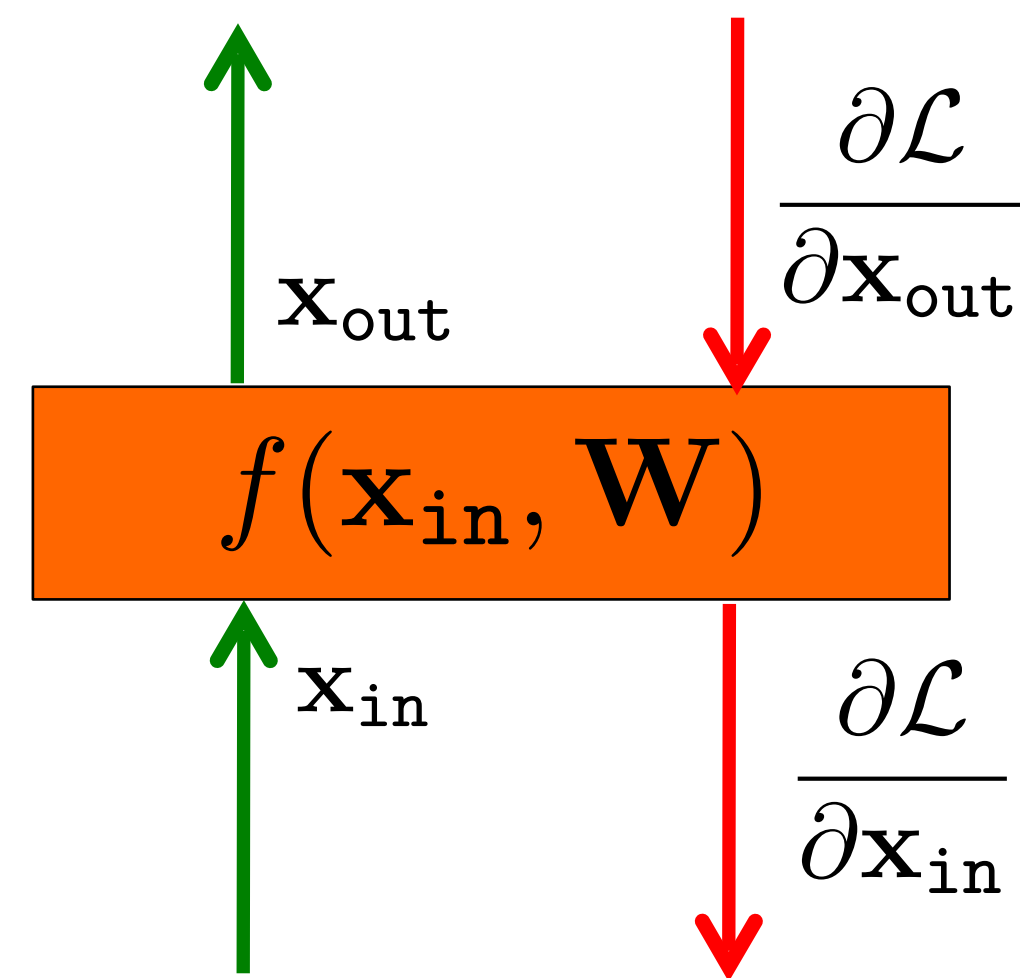


- Forward propagation: $\mathbf{x}_{\text{out}} = f(\mathbf{x}_{\text{in}}, \mathbf{W}) = \mathbf{W}\mathbf{x}_{\text{in}}$
- Backprop to input:



Now let's see how we use the set of outputs to compute the weights update equation (backprop to the weights).

Linear Module



- Forward propagation: $\mathbf{x}_{\text{out}} = f(\mathbf{x}_{\text{in}}, \mathbf{W}) = \mathbf{W}\mathbf{x}_{\text{in}}$
- Backprop to weights:

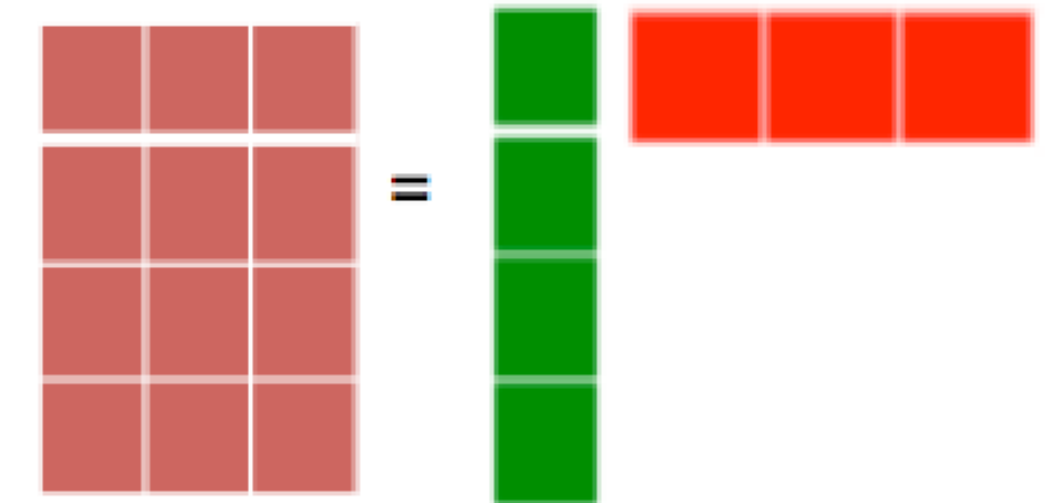
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}} \cdot \frac{\partial f(\mathbf{x}_{\text{in}}, \mathbf{W})}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}} \cdot \frac{\partial \mathbf{x}_{\text{out}}}{\partial \mathbf{W}}$$

If we look at how the parameter W_{ij} changes the cost, only the i component of the output will change, therefore:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{ij}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot \frac{\partial \mathbf{x}_{\text{out}_i}}{\partial \mathbf{W}_{ij}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot \mathbf{x}_{\text{in}_j}$$

\uparrow
 $\frac{\partial \mathbf{x}_{\text{out}_i}}{\partial \mathbf{W}_{ij}} = \mathbf{x}_{\text{in}_j}$

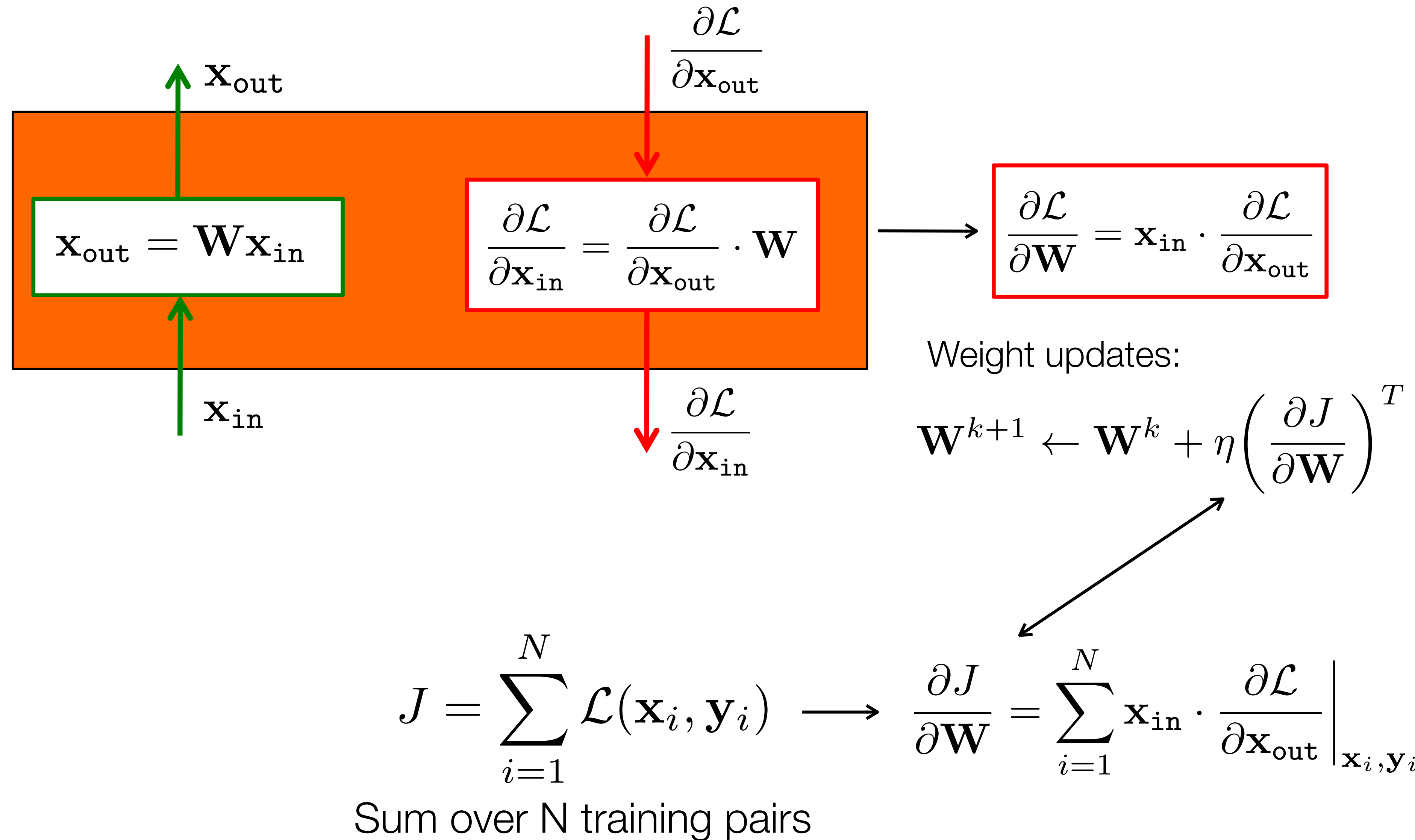
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \mathbf{x}_{\text{in}} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}}}$$



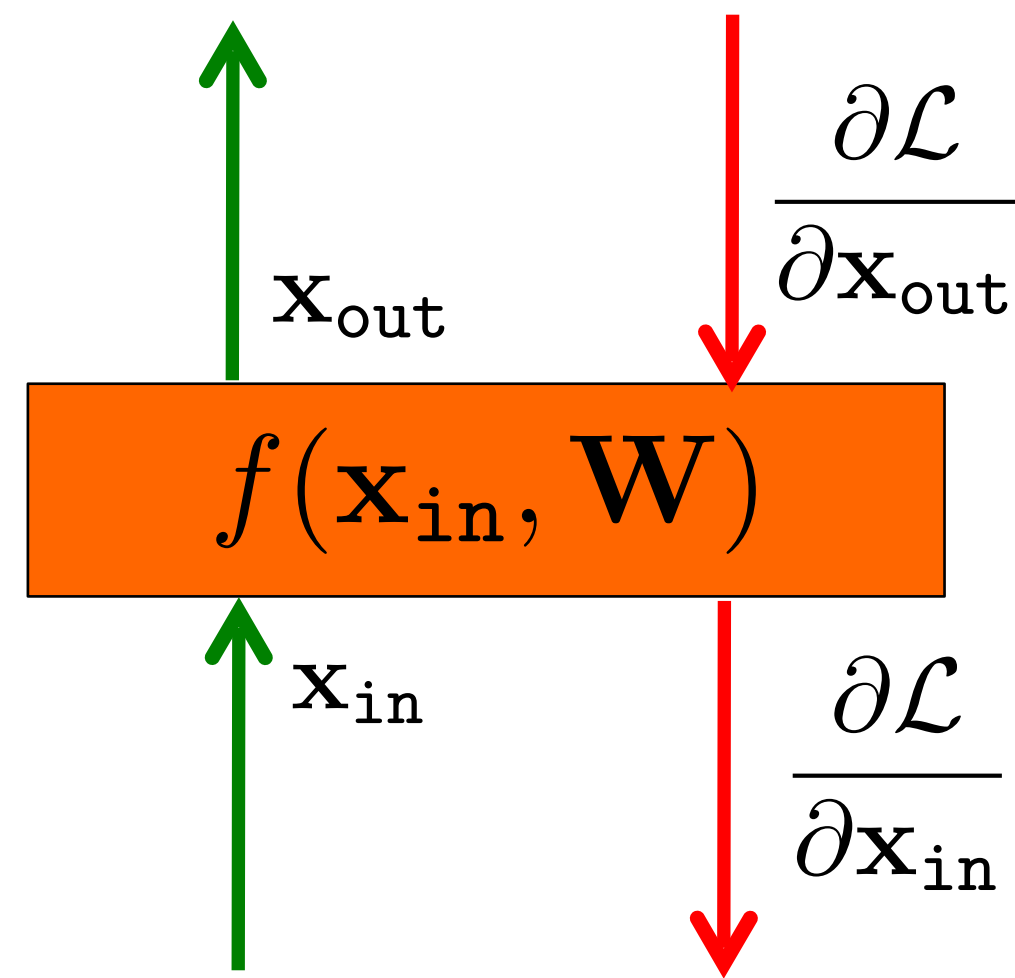
And now we can update the weights (by summing over all the training examples):

$$\mathbf{W}^{k+1} \leftarrow \mathbf{W}^k + \eta \left(\frac{\partial J}{\partial \mathbf{W}} \right)^T \quad (\text{sum over all training examples to get } J)$$

Linear Module



Pointwise function



- Forward propagation:

$$\mathbf{x}_{\text{out}_i} = h(\mathbf{x}_{\text{in}_i} + b_i)$$

h = an arbitrary function, b_i is a bias term.

- Backprop to input: $\frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{in}_i}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot \frac{\partial \mathbf{x}_{\text{out}_i}}{\partial \mathbf{x}_{\text{in}_i}} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot h'(\mathbf{x}_{\text{in}_i} + b_i)$
- Backprop to bias: $\frac{\partial \mathcal{L}}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot \frac{\partial \mathbf{x}_{\text{out}_i}}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_{\text{out}_i}} \cdot h'(\mathbf{x}_{\text{in}_i} + b_i)$

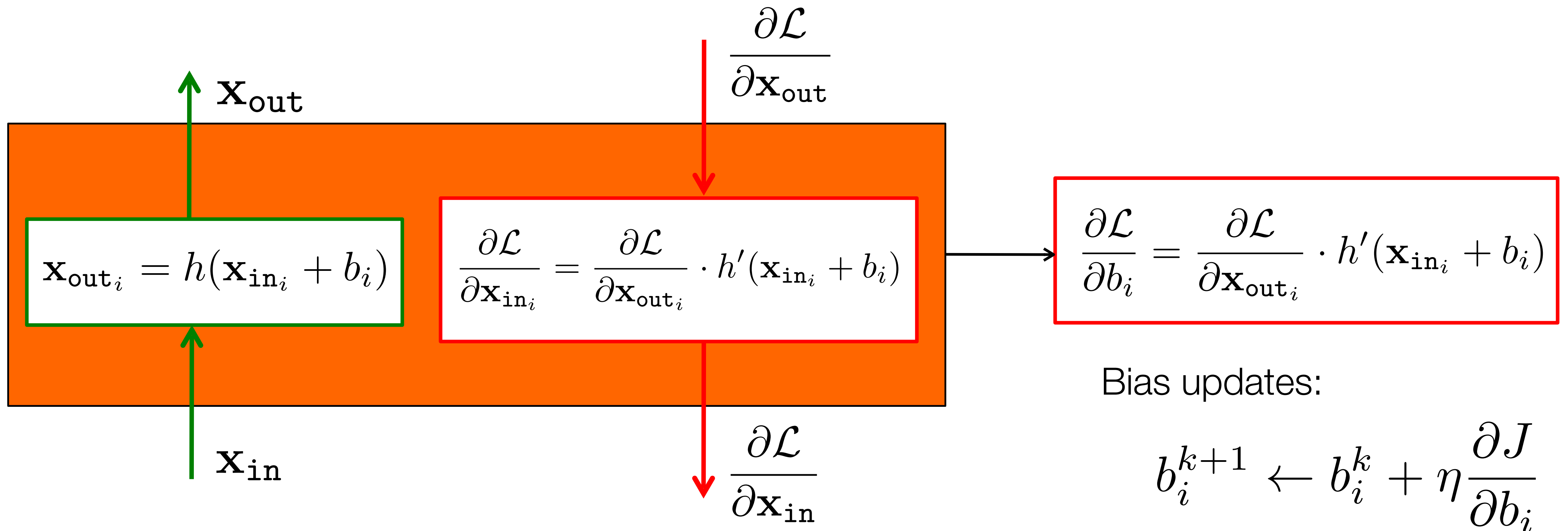
We use this last expression to update the bias.

Some useful derivatives:

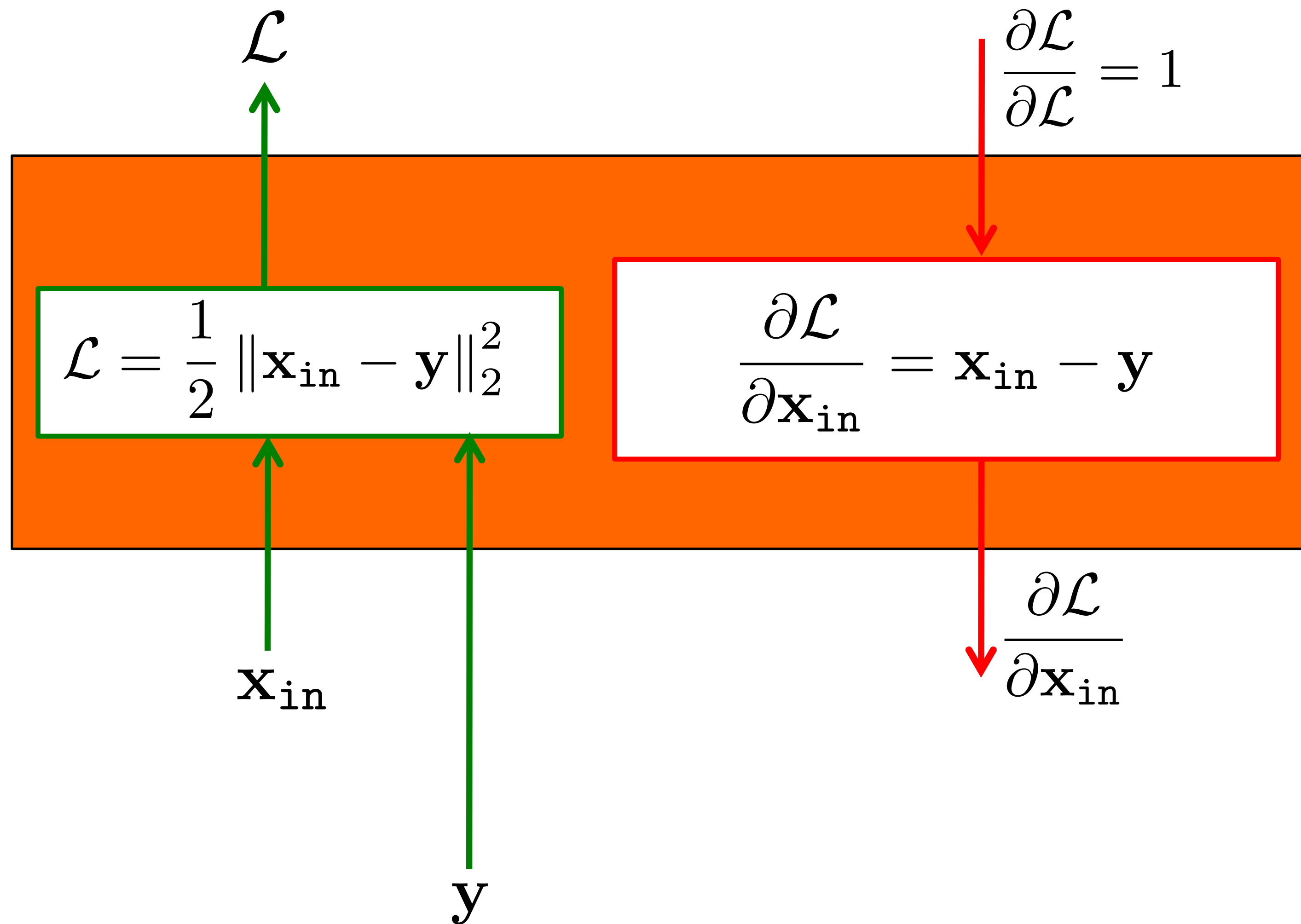
For hyperbolic tangent: $\tanh'(x) = 1 - \tanh^2(x)$

For ReLU: $h(x) = \max(0, x)$, $h'(x) = \mathbb{1}(x \geq 0)$

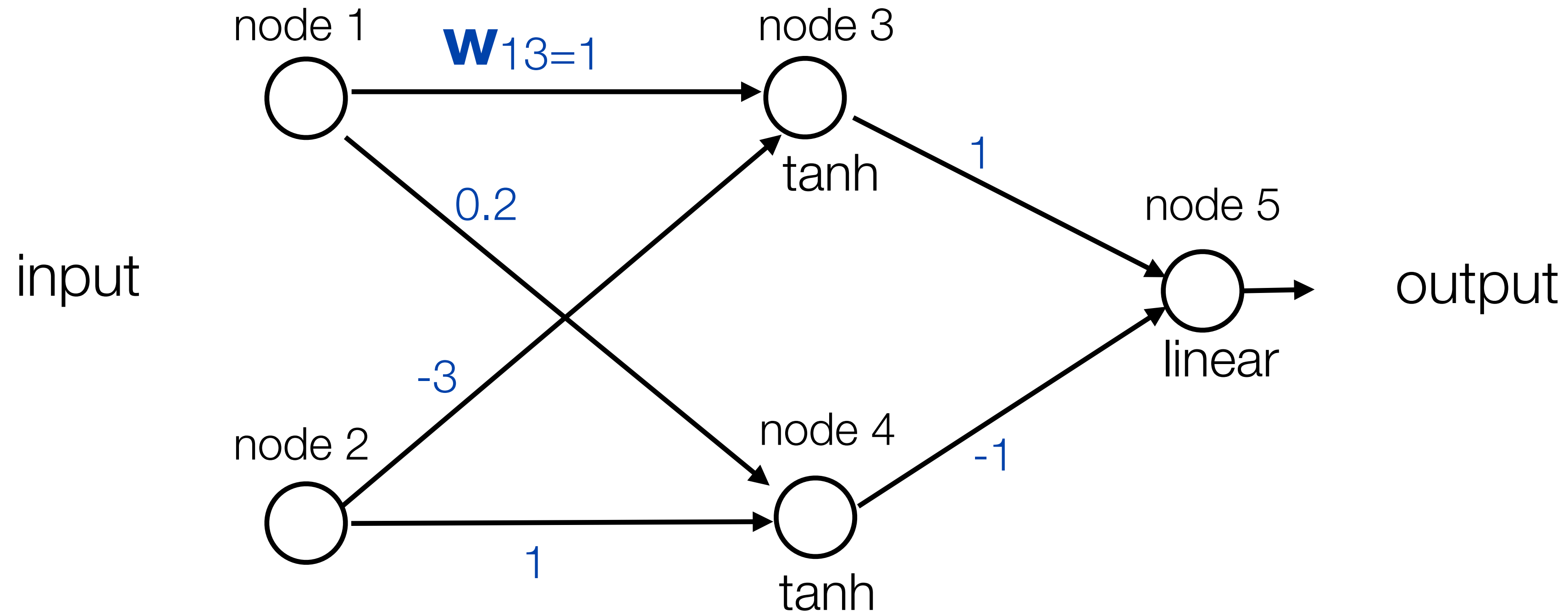
Pointwise function



Euclidean cost module



Backpropagation example



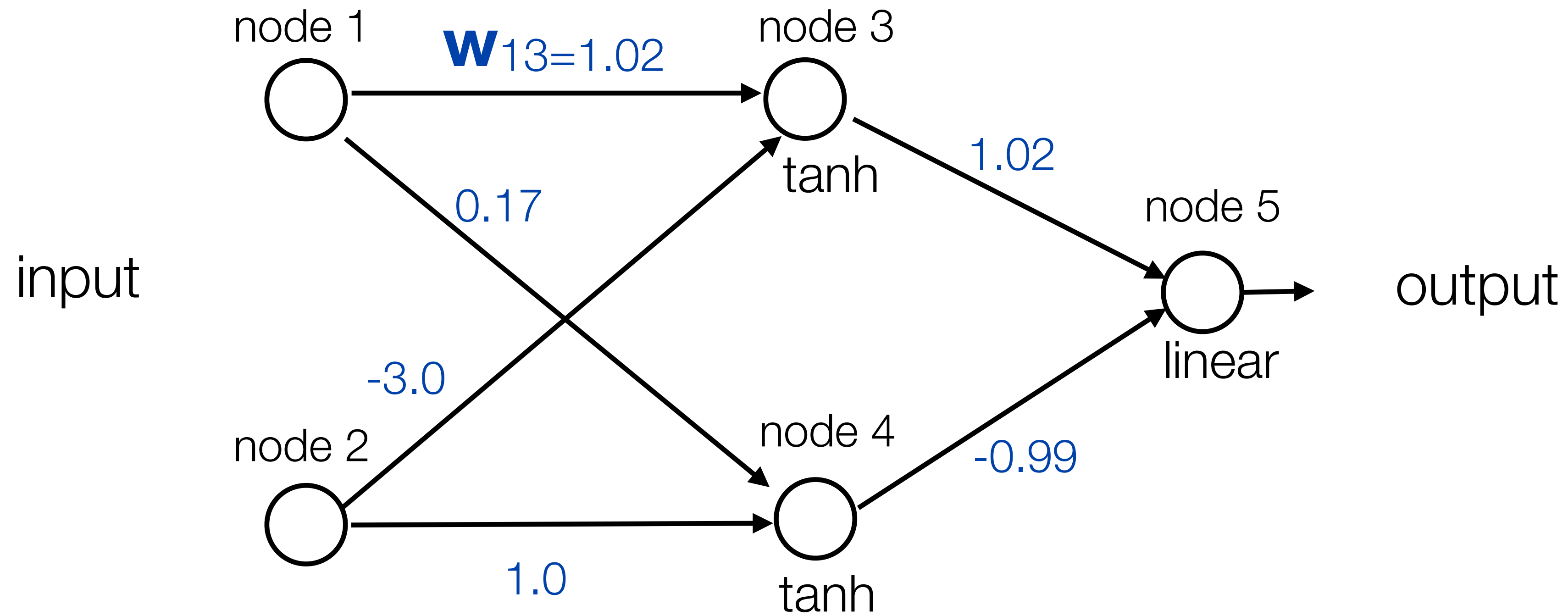
Learning rate $\eta = -0.2$ (because we used positive increments)

Euclidean loss

Training data: input		desired output
node 1	node 2	node 5
1.0	0.1	0.5

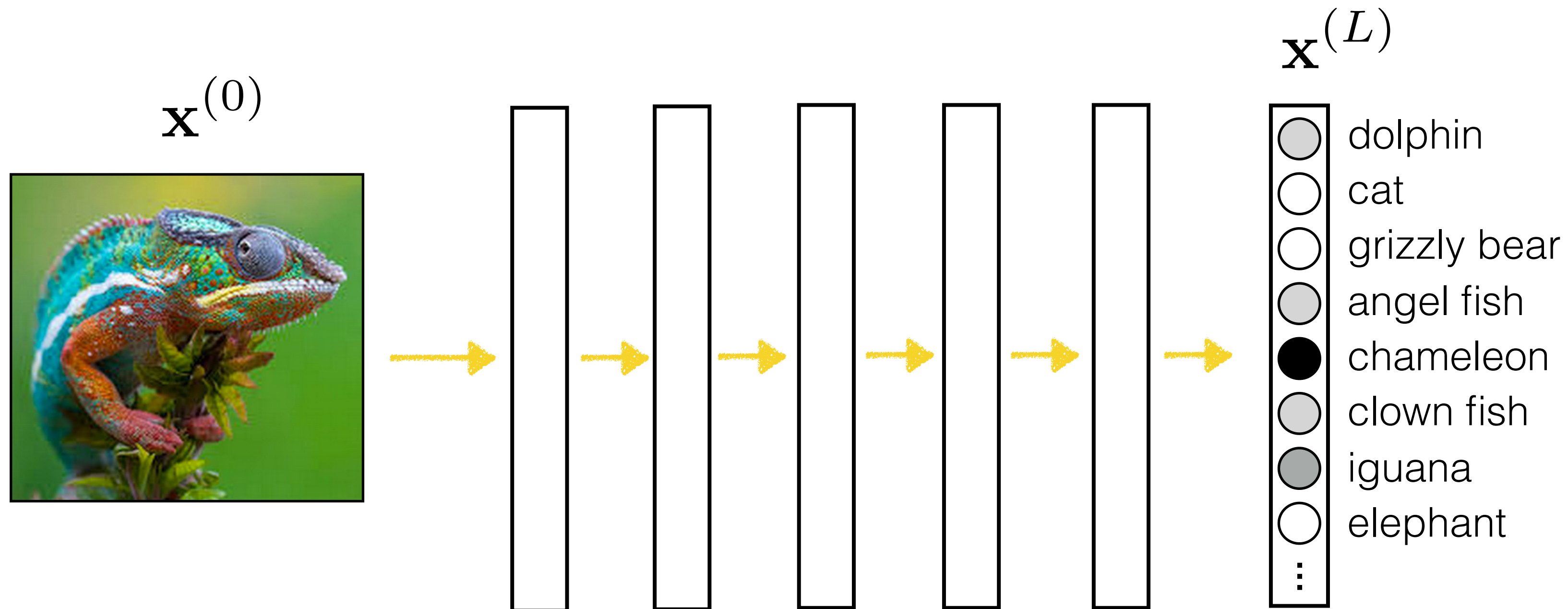
Exercise: run one iteration of back propagation

Backpropagation example

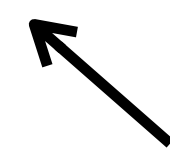


After one iteration (rounding to two digits)

Unit visualization via backprop



$$\frac{\partial x_j^{(L)}}{\partial \mathbf{x}^{(0)}} = \frac{\partial x_j^{(L)}}{\partial \mathbf{x}^{(L-1)}} \cdot \frac{\partial \mathbf{x}^{(L-1)}}{\partial \mathbf{x}^{(L-2)}} \cdots \frac{\partial \mathbf{x}^{(2)}}{\partial \mathbf{x}^{(1)}} \cdot \frac{\partial \mathbf{x}^{(1)}}{\partial \mathbf{x}^{(0)}}$$



How much the “chameleon” score is increased or decreased by changing the image pixels.

Unit visualization via backprop

$$\arg \max_{\mathbf{x}^{(0)}} x_j^{(L)}$$

$$\mathbf{x}^{(0)^{k+1}} \leftarrow \mathbf{x}^{(0)^k} + \eta \frac{\partial x_j^{(L)}}{\partial \mathbf{x}^{(0)}}$$

Unit visualization via backprop

Make an image that maximizes the “cat”
output neuron:

$$\arg \max_{\mathbf{x}^{(0)}} x_j^{(L)} + \lambda R(\mathbf{x}^{(0)})$$

$$\mathbf{x}^{(0)^{k+1}} \leftarrow \mathbf{x}^{(0)^k} + \eta \frac{\partial (x_j^{(L)} + \lambda R(\mathbf{x}^{(0)}))}{\partial \mathbf{x}^{(0)}}$$



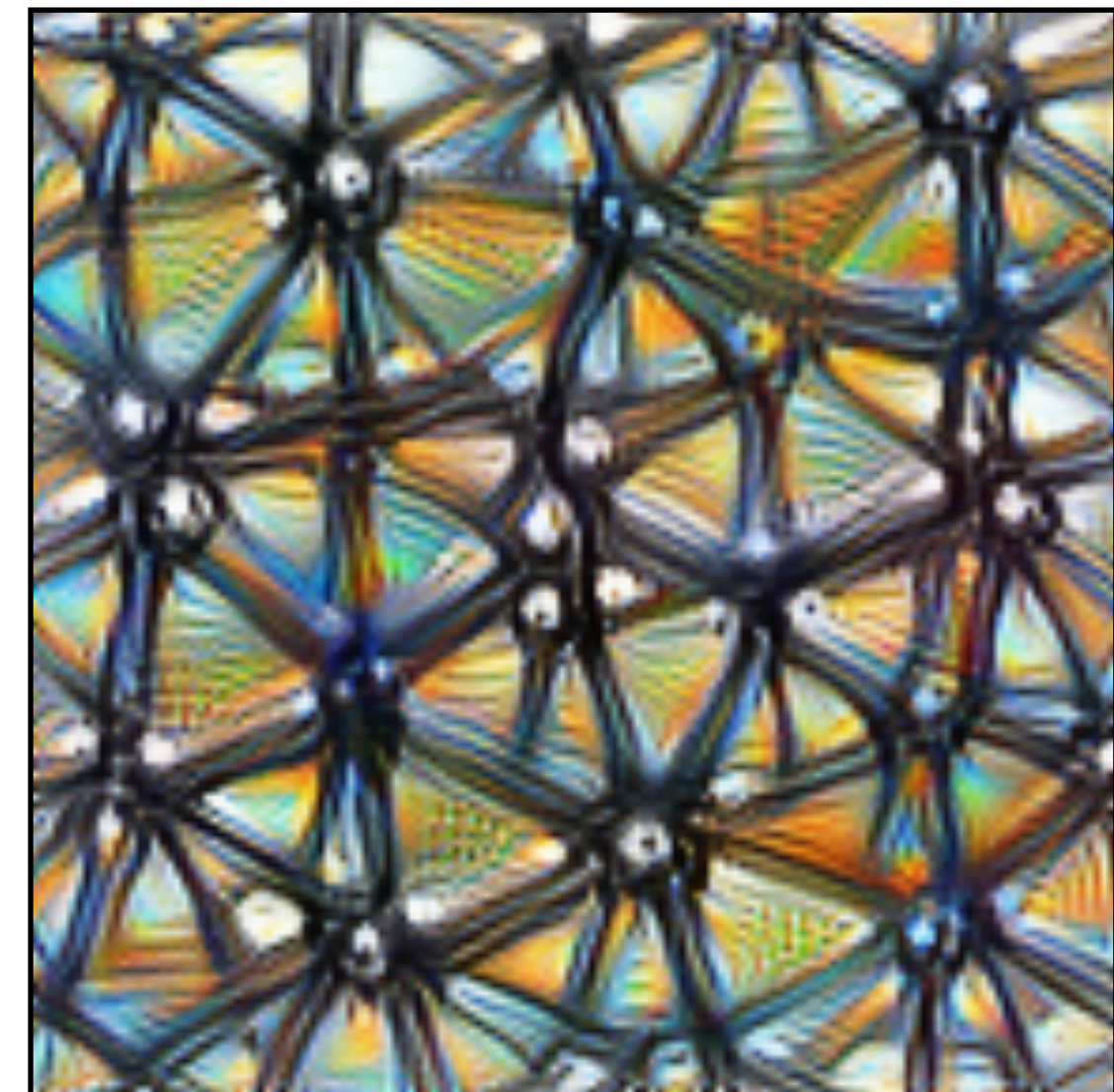
[<https://distill.pub/2017/feature-visualization/>]

Unit visualization via backprop

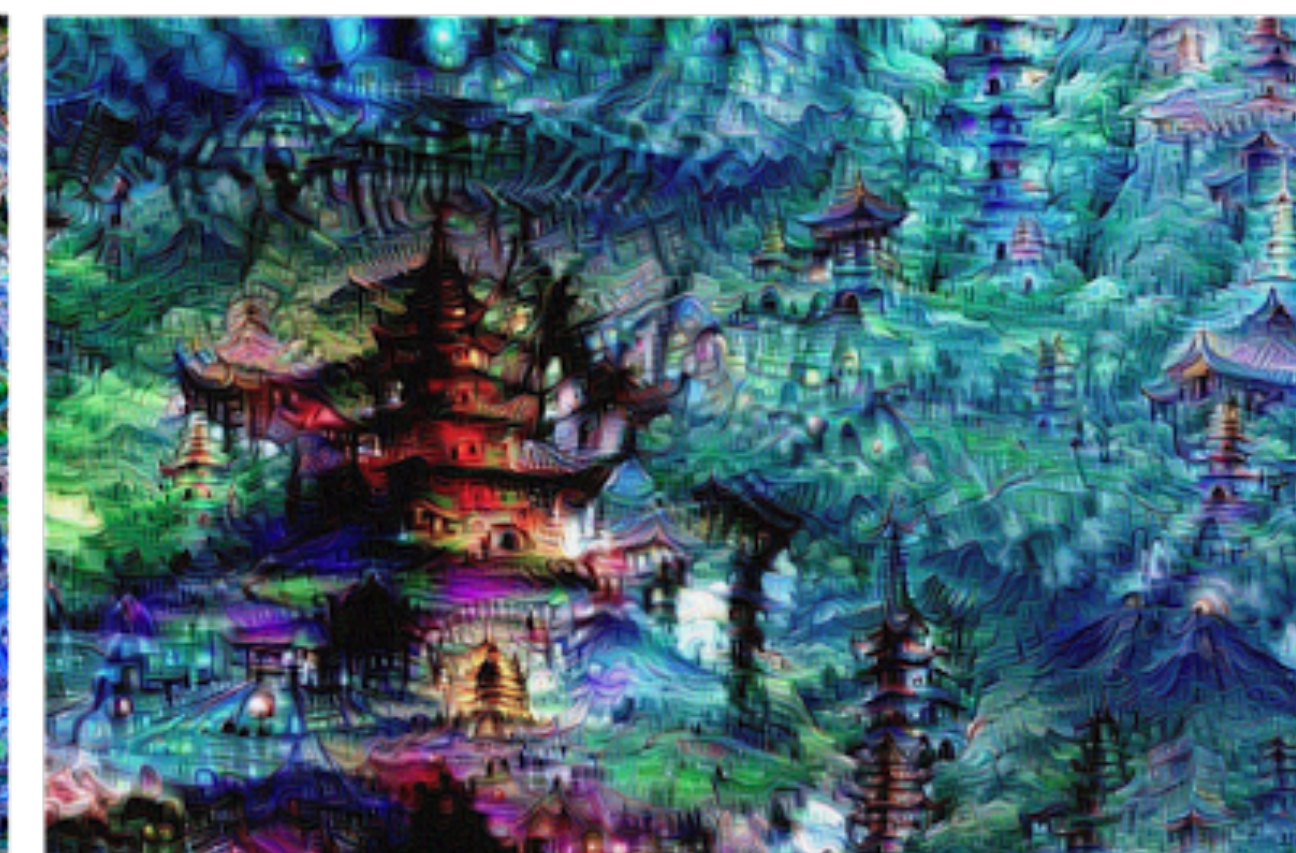
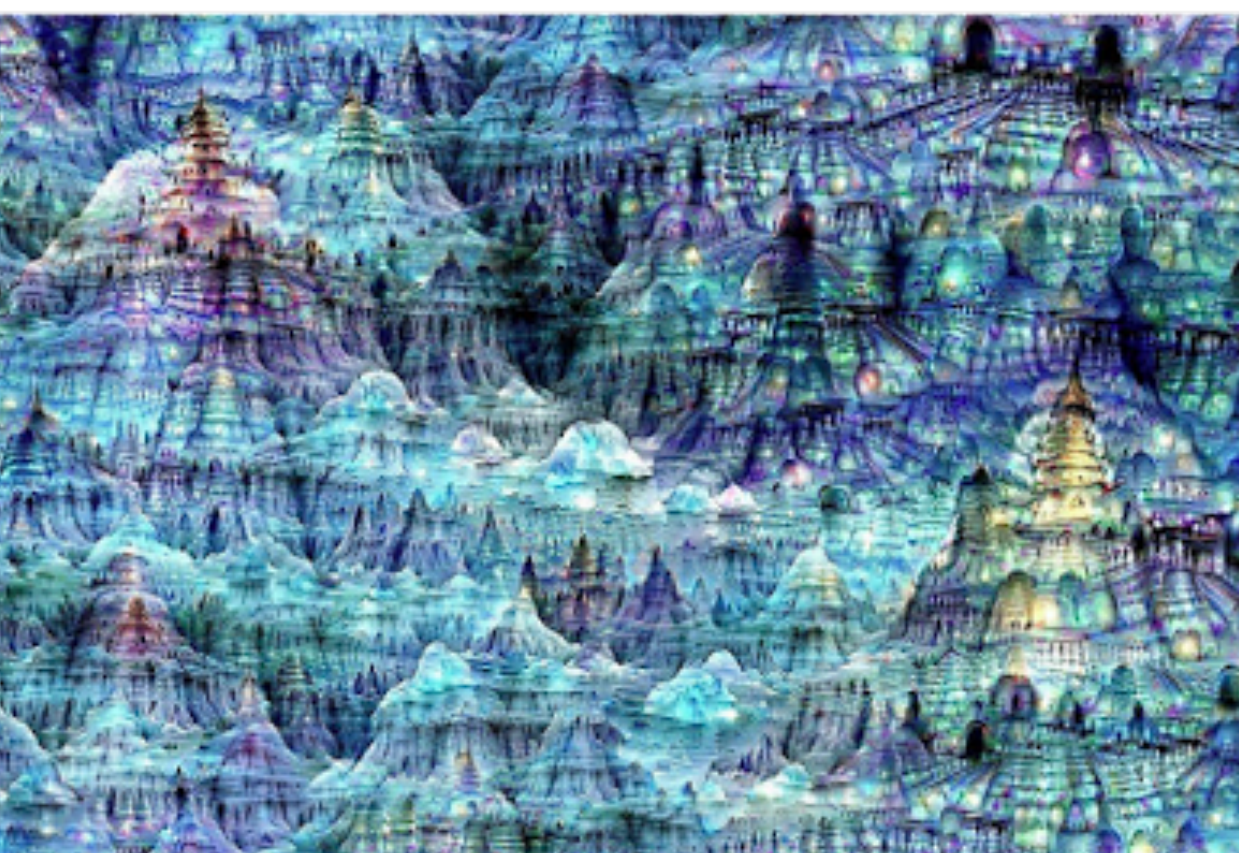
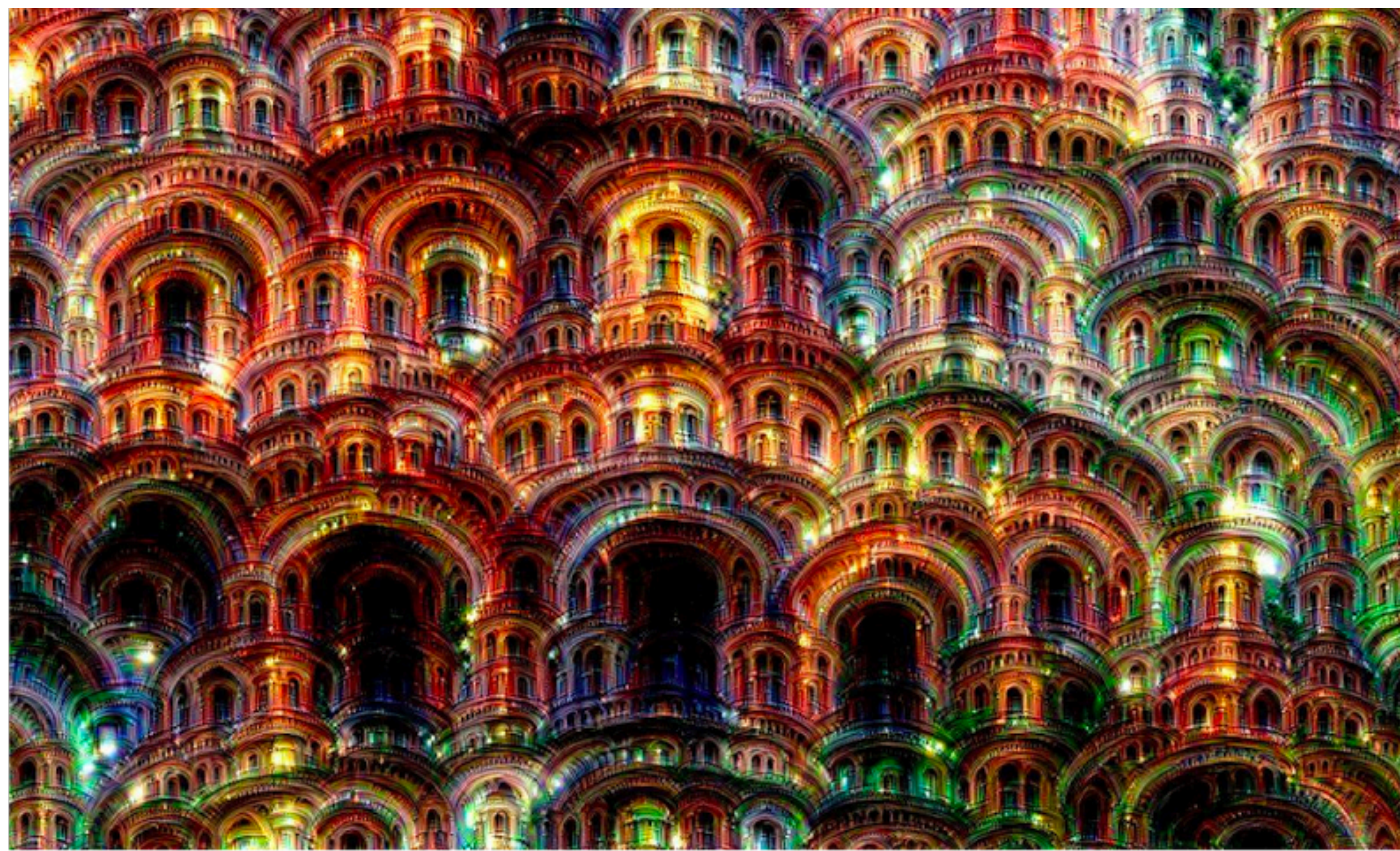
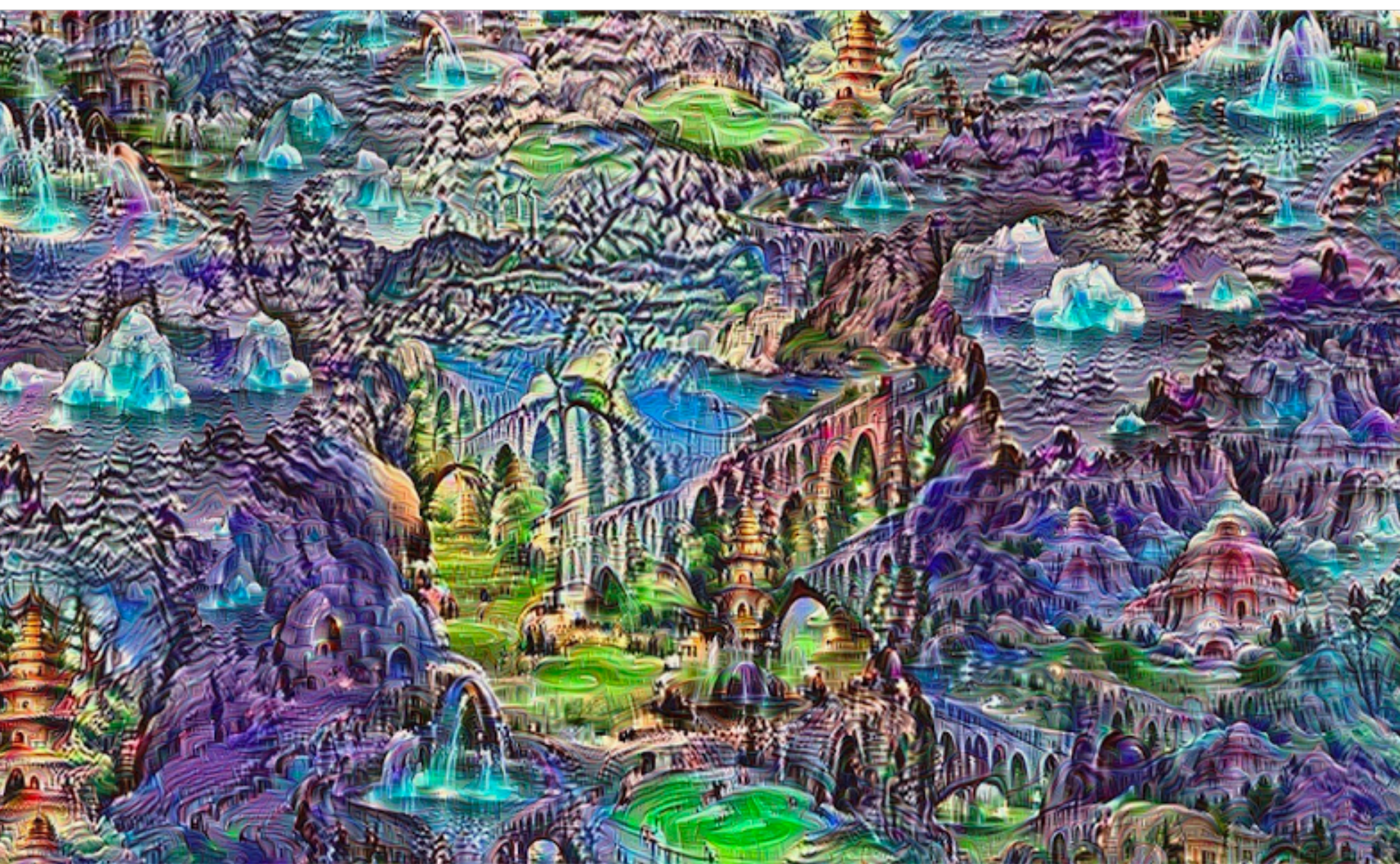
Make an image that maximizes the value of a random neuron in the middle of the network:

$$\arg \max_{\mathbf{x}^{(0)}} x_j^{(L)} + \lambda R(\mathbf{x}^{(0)})$$

$$\mathbf{x}^{(0)^{k+1}} \leftarrow \mathbf{x}^{(0)^k} + \eta \frac{\partial (x_j^{(L)} + \lambda R(\mathbf{x}^{(0)}))}{\partial \mathbf{x}^{(0)}}$$



[<https://distill.pub/2017/feature-visualization/>]



“Deep dream” [<https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>]

Differentiable programming

Deep nets are popular for a few reasons:

1. High capacity
2. Easy to optimize (differentiable)
3. Compositional “block based programming”

An emerging term for general models with these properties is **differentiable programming**.



Yann LeCun

January 5 · 🌐

OK, Deep Learning has outlived its usefulness as a buzz-phrase. Deep Learning est mort. Vive Differentiable Programming!



Thomas G. Dietterich

@tdietterich

Following

DL is essentially a new style of programming--"differentiable programming"--and the field is trying to work out the reusable constructs in this style. We have some: convolution, pooling, LSTM, GAN, VAE, memory units, routing units, etc. 8/

8:02 AM - 4 Jan 2018

65 Retweets 194 Likes



6

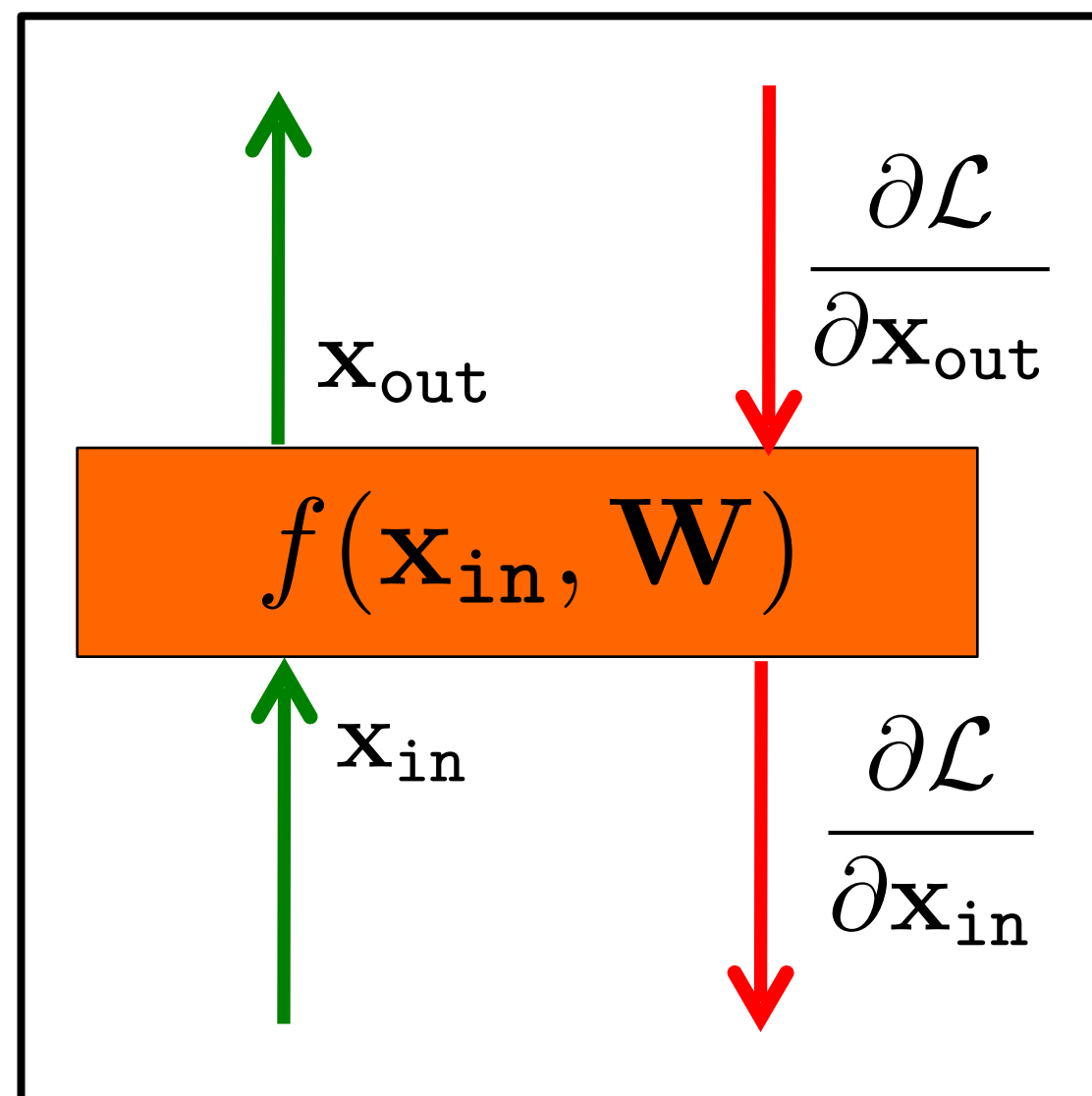
65

194

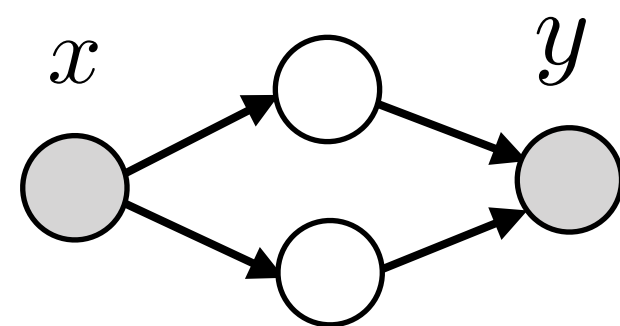
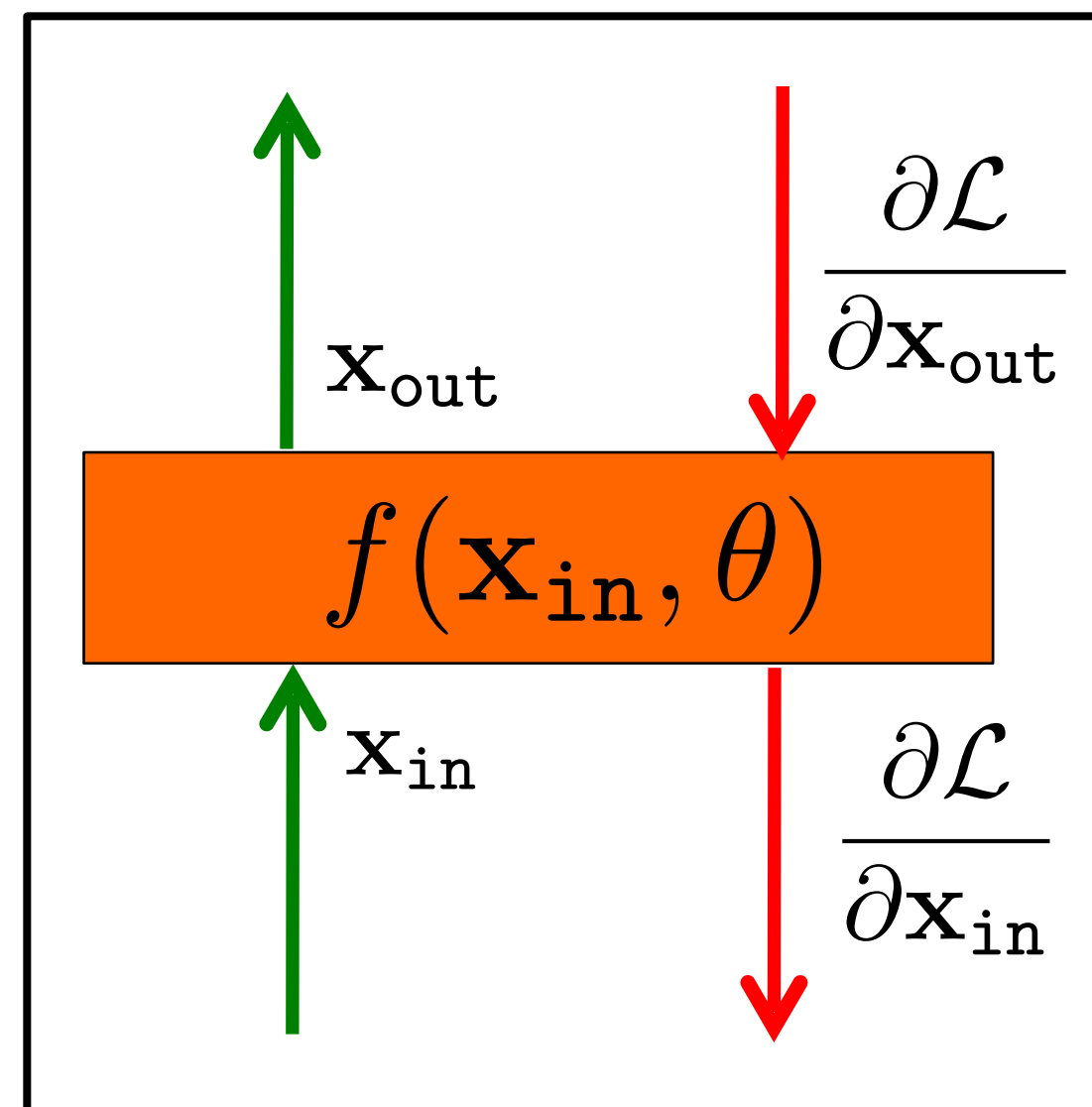


Differentiable programming

Deep learning

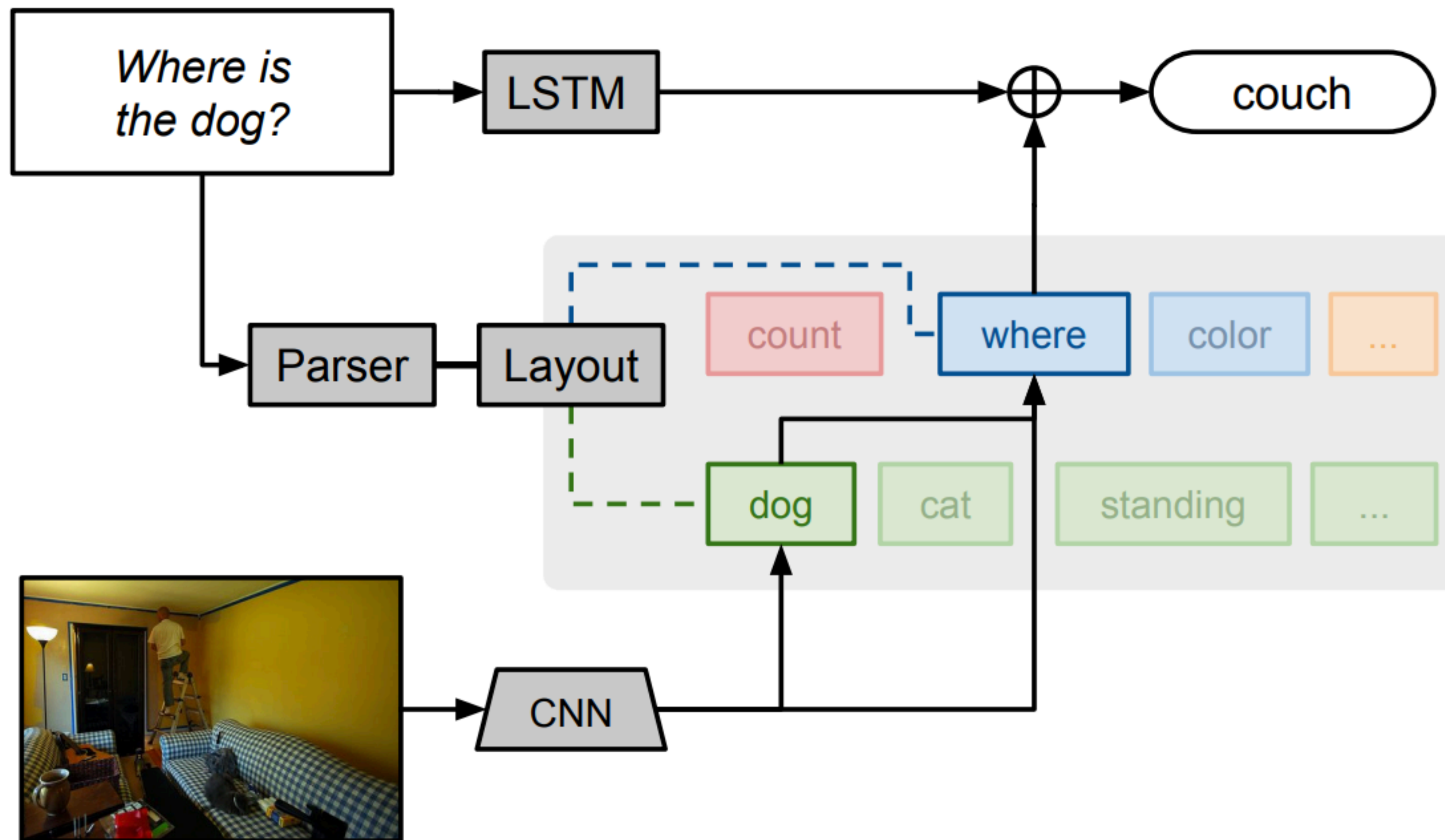


Differentiable programming



```
1 for i, data in enumerate(dataset):
2     iter_start_time = time.time()
3     if total_steps % opt.print_freq == 0:
4         t_data = iter_start_time - iter_data_time
5         visualizer.reset()
6         total_steps += opt.batch_size
7         epoch_iter += opt.batch_size
8         model.set_input(data)
9         model.optimize_parameters()
```

Differentiable programming



[Figure from "Neural Module Networks", Andreas et al. 2017]

Next up:

1. Convolutional neural networks (CNNs)