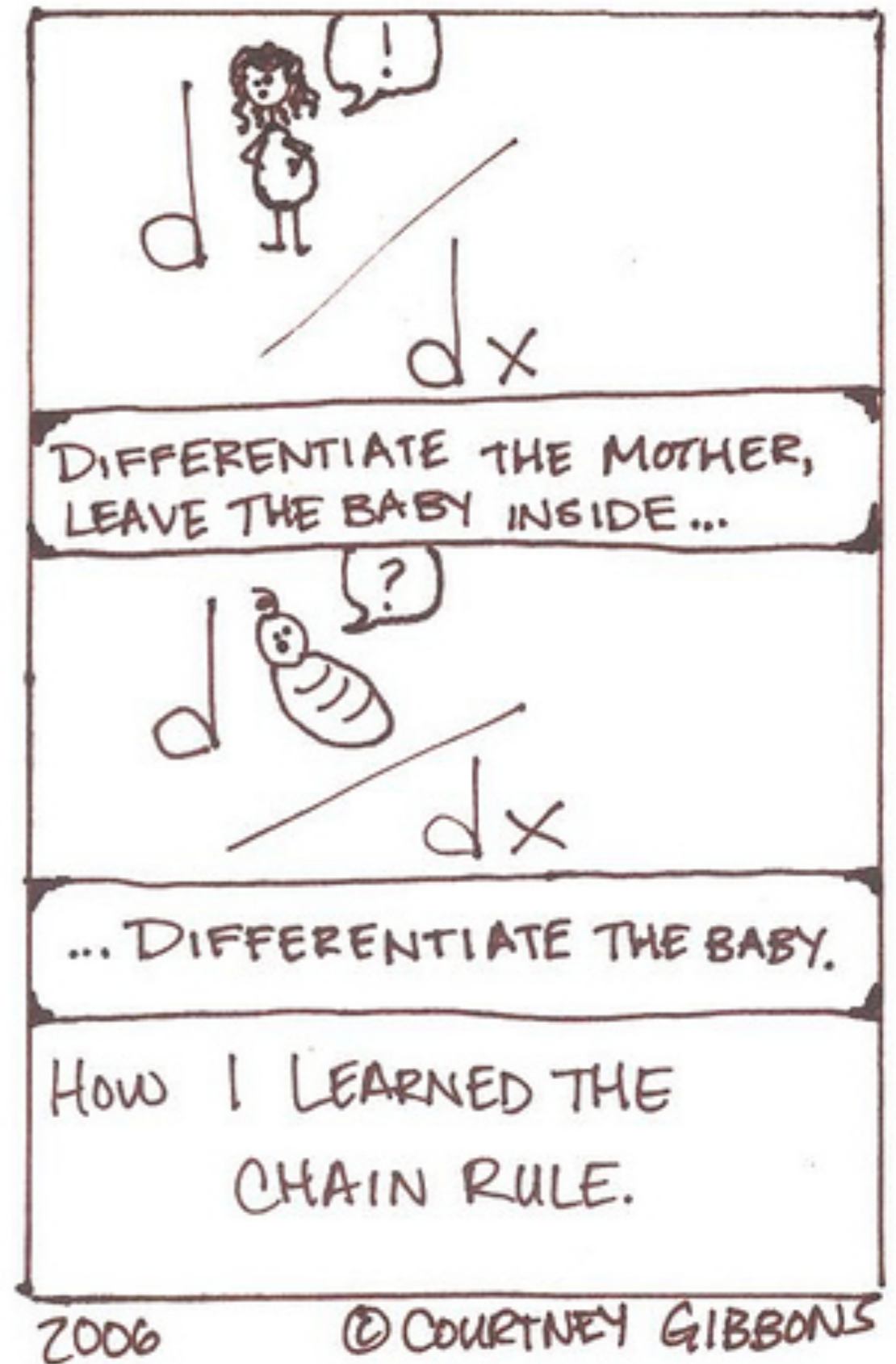
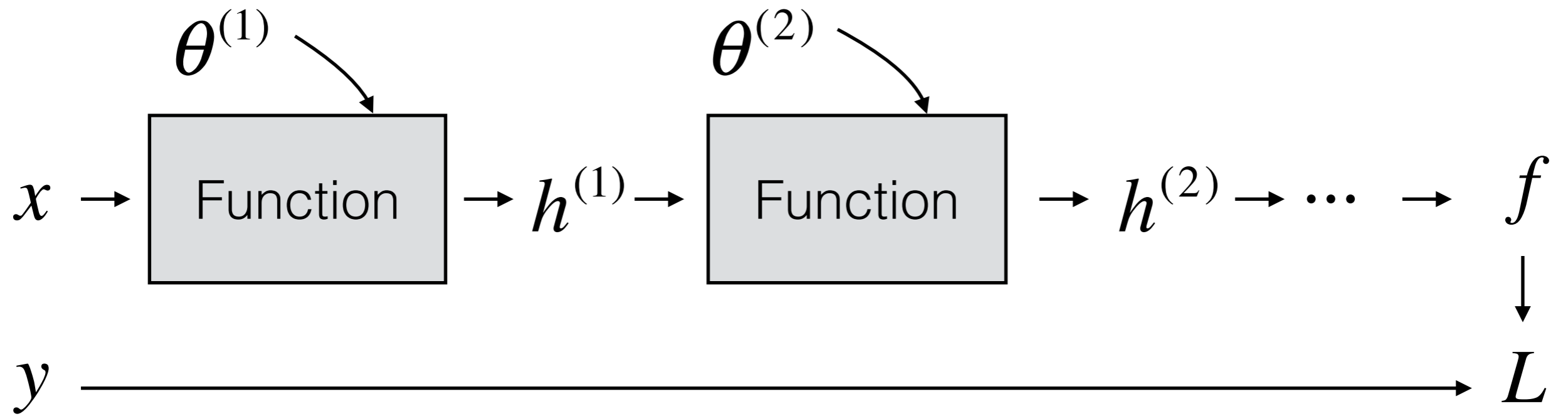


Convolutional Neural Networks

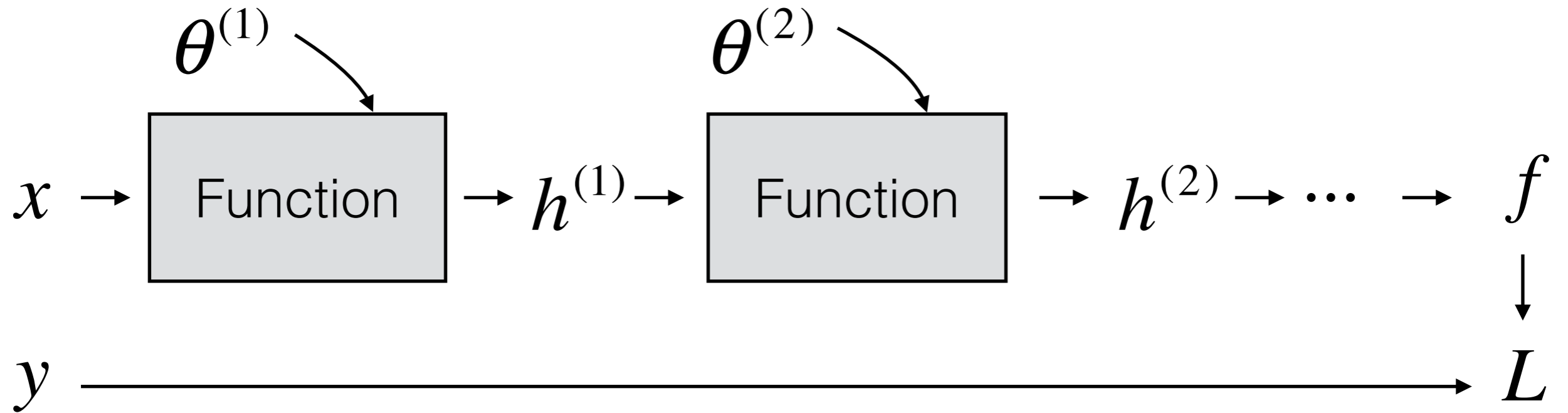
CS 4670
Sean Bell



Review: Setup

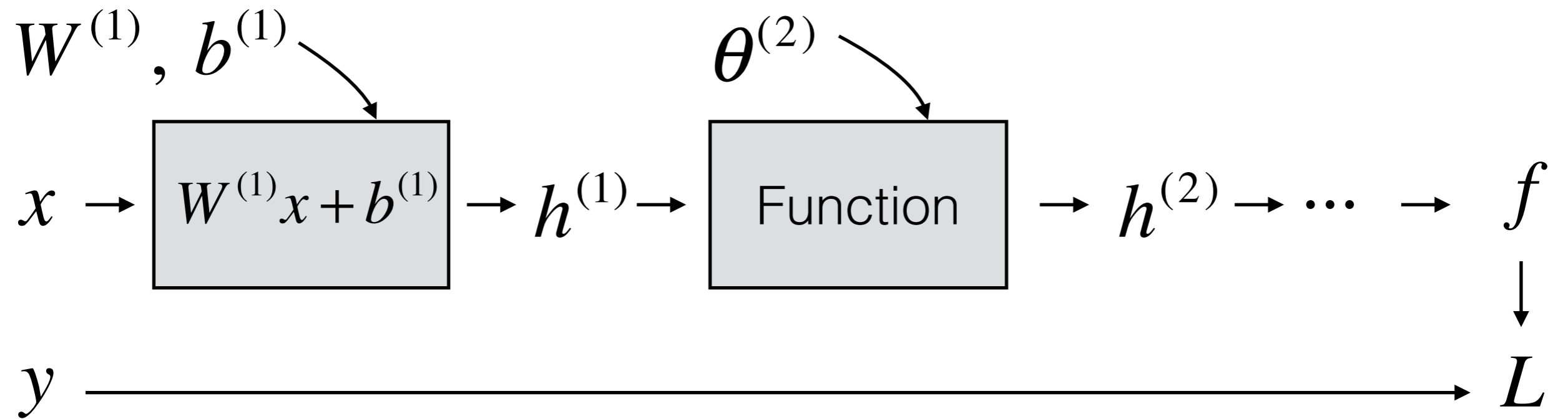


Review: Setup



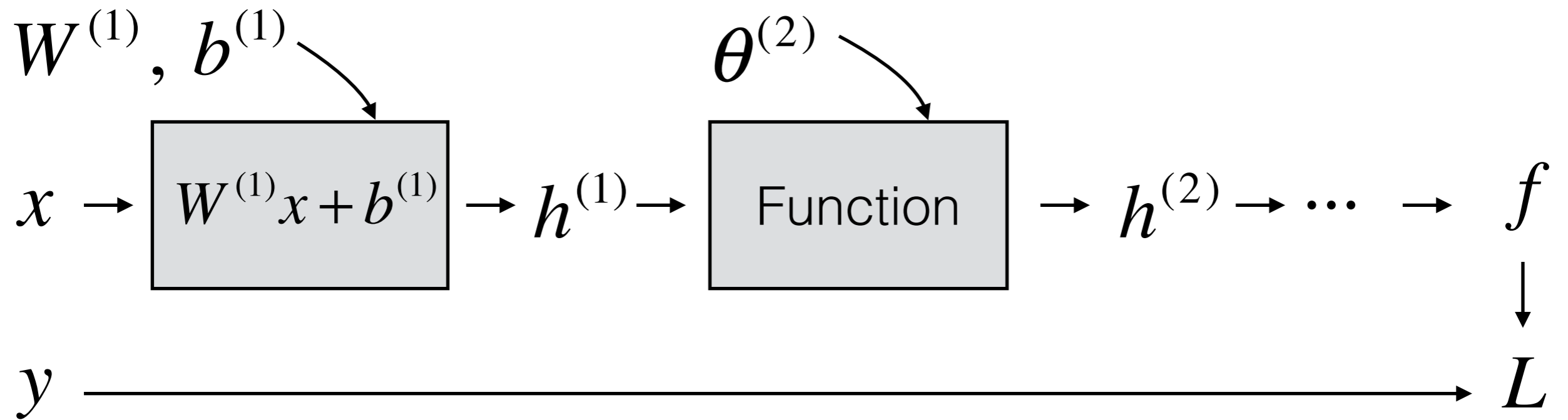
- **Goal:** Find a value for parameters $(\theta^{(1)}, \theta^{(2)}, \dots)$, so that the loss (L) is small

Review: Setup

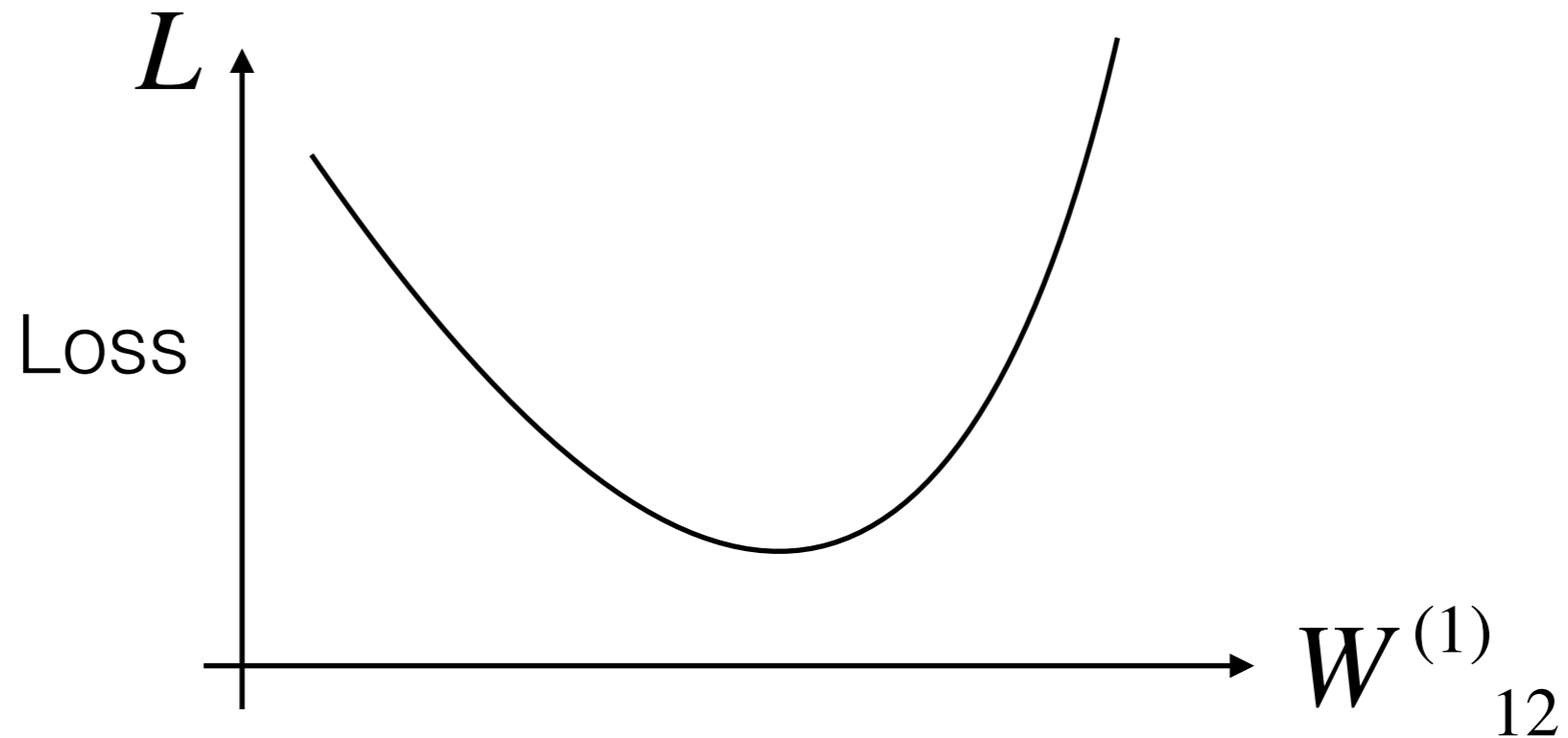


**Toy
Example:**

Review: Setup

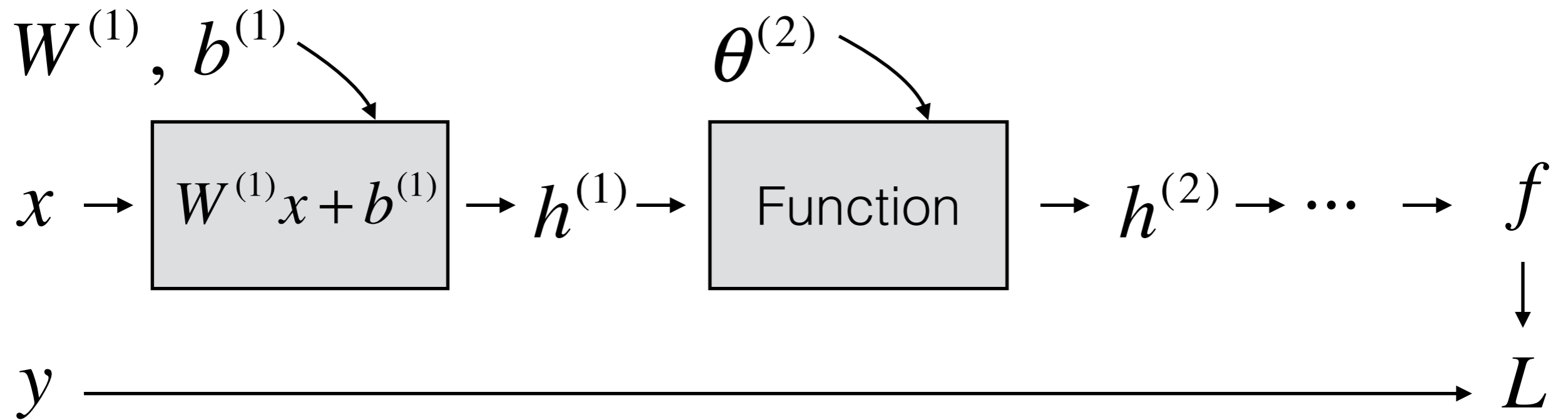


Toy Example:

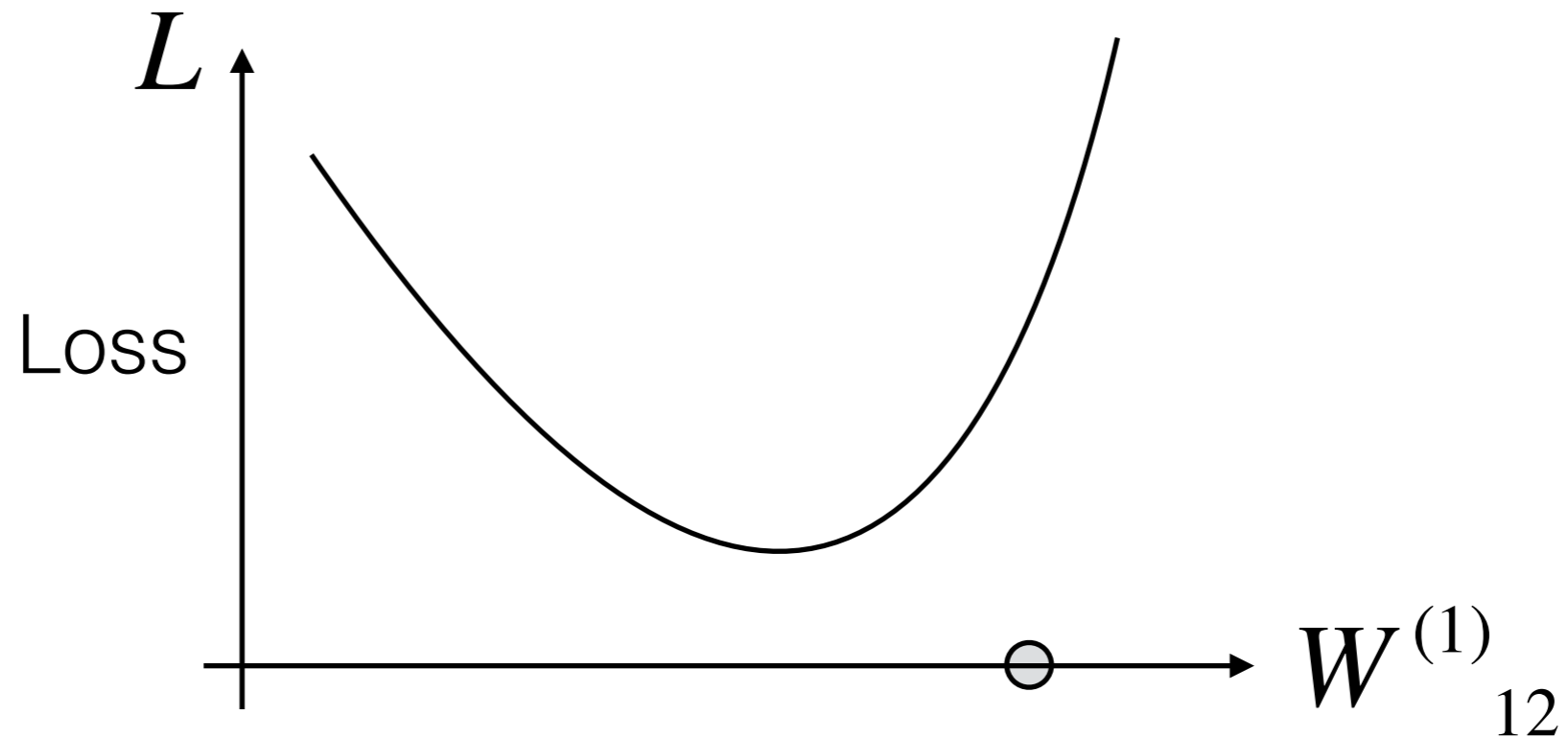


A weight somewhere in the network

Review: Setup

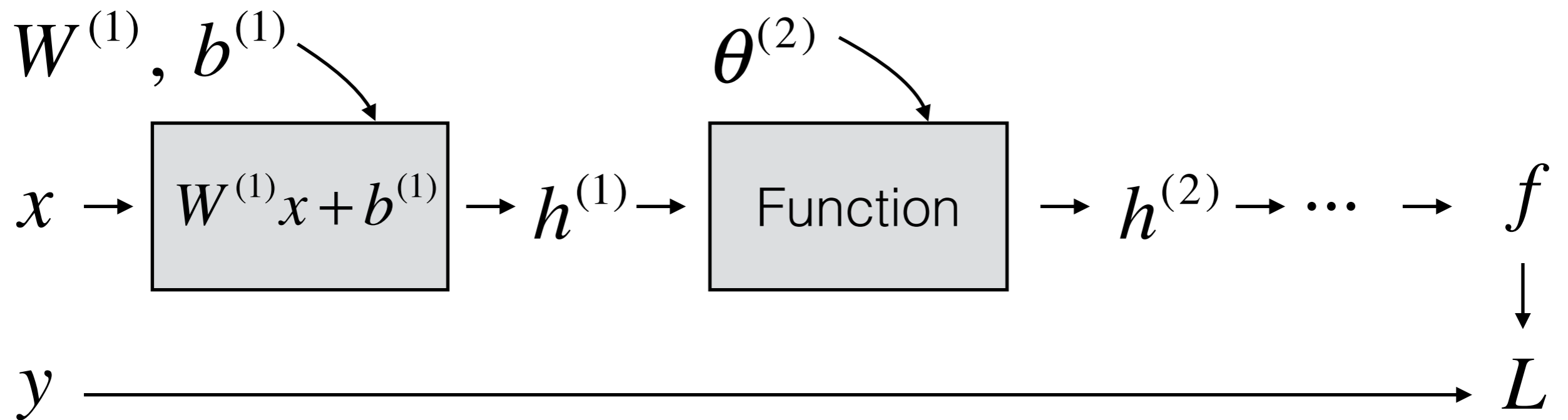


Toy Example:

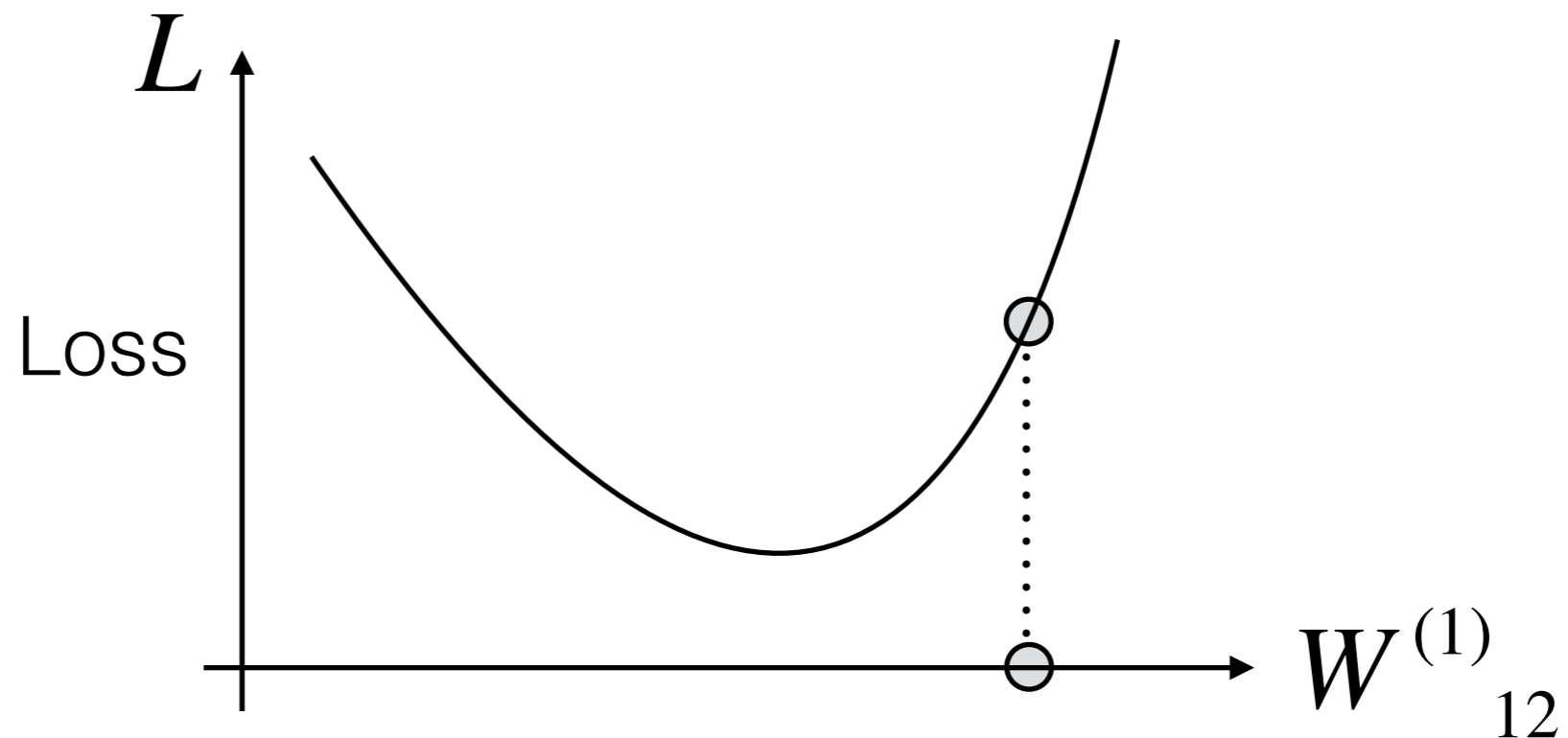


A weight somewhere in the network

Review: Setup

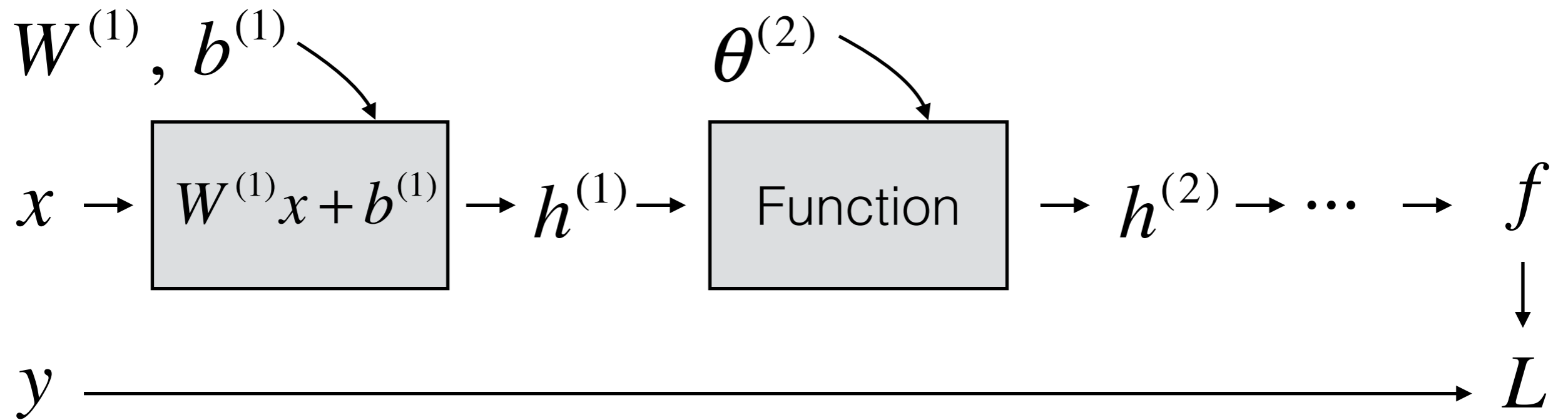


Toy Example:

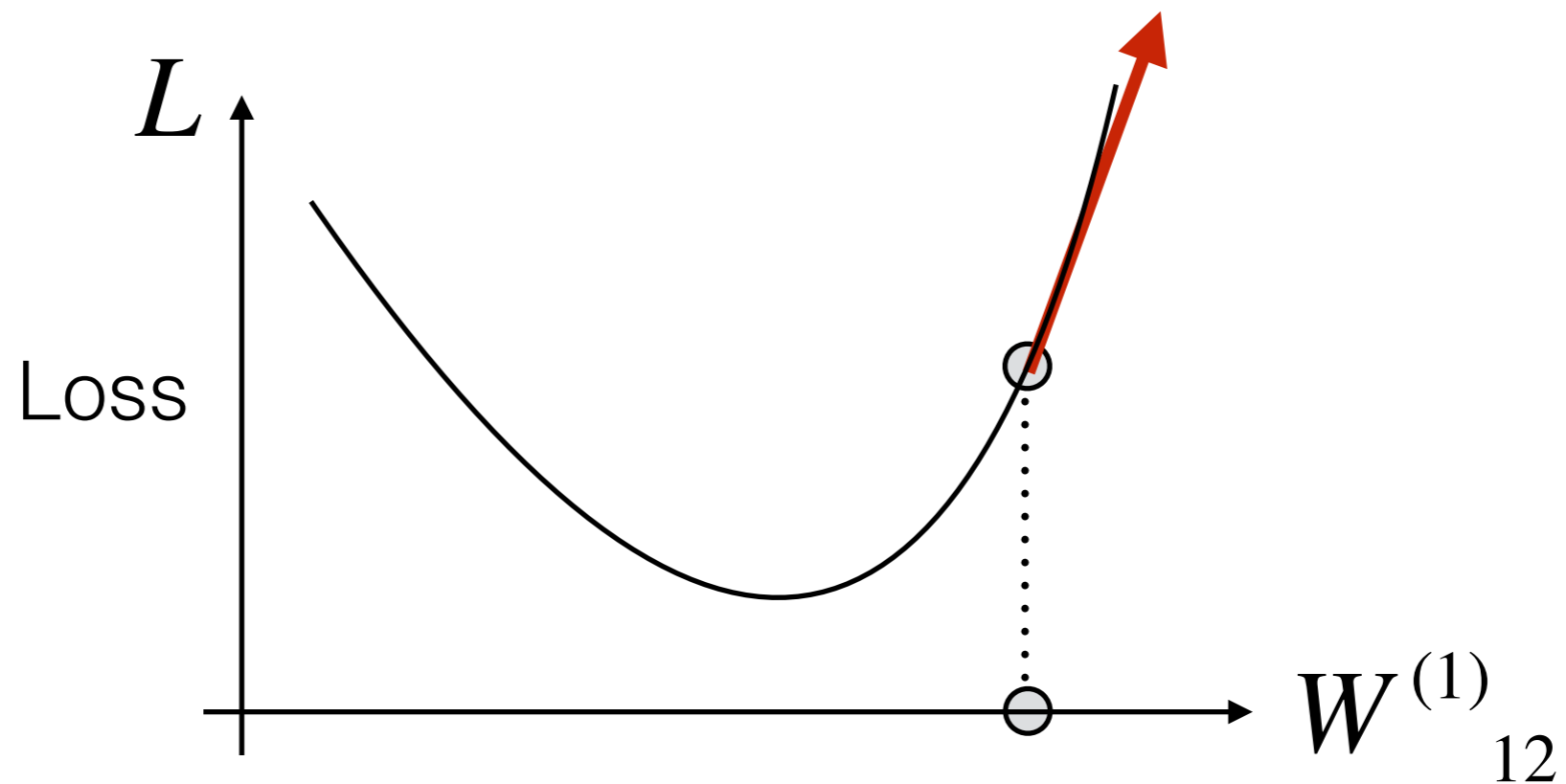


A weight somewhere in the network

Review: Setup

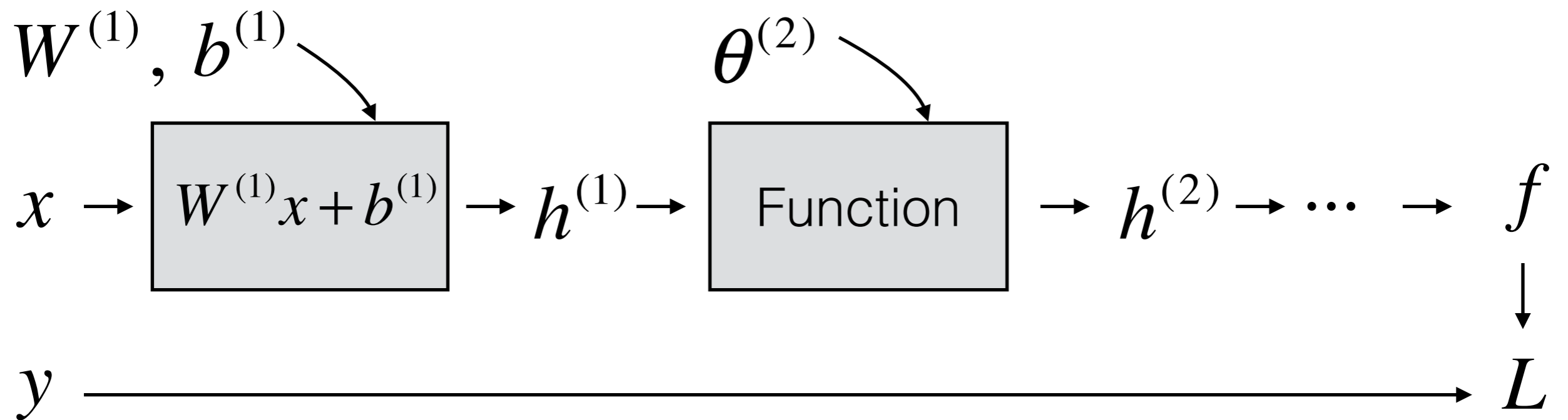


Toy Example:

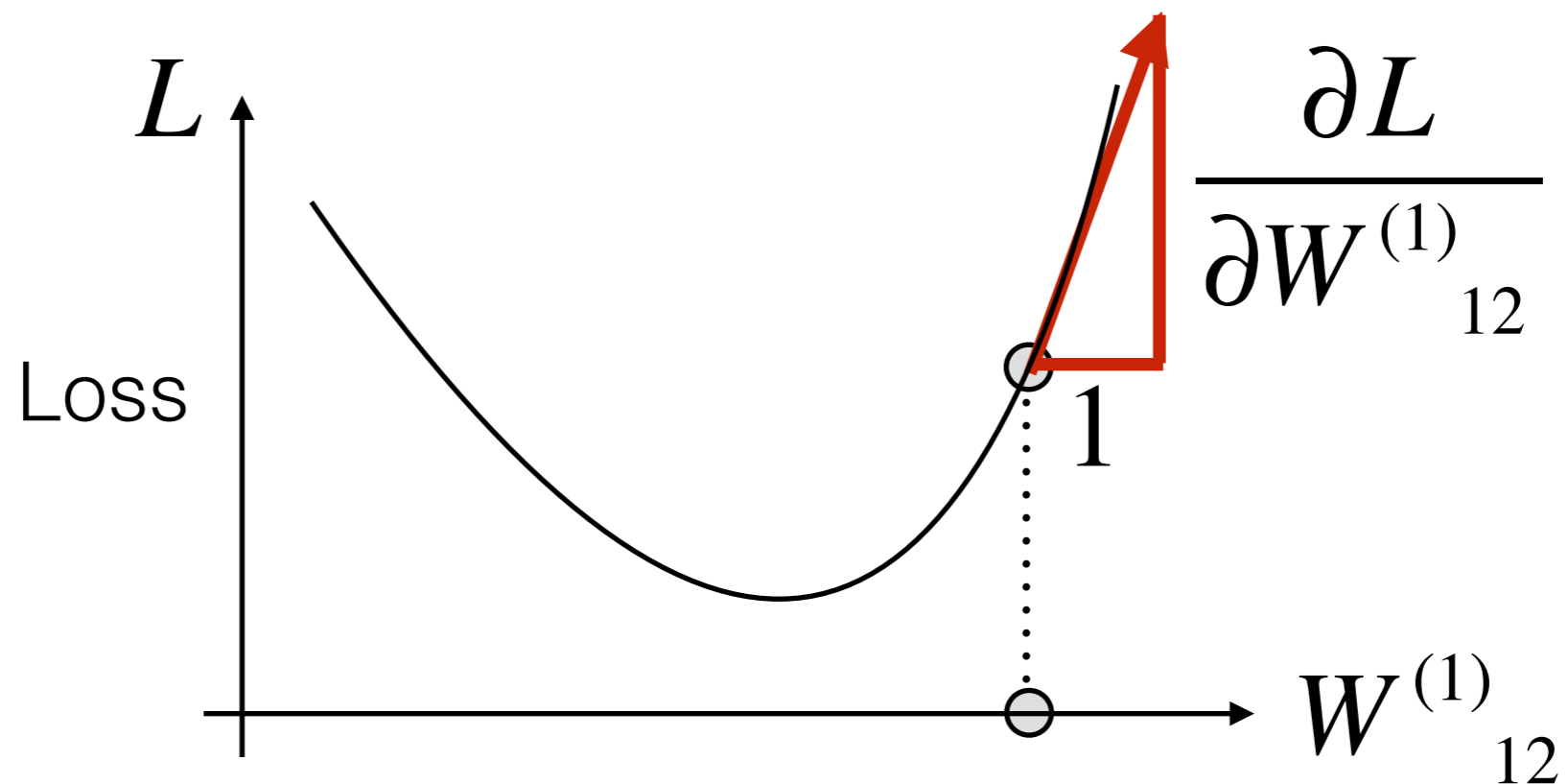


A weight somewhere in the network

Review: Setup

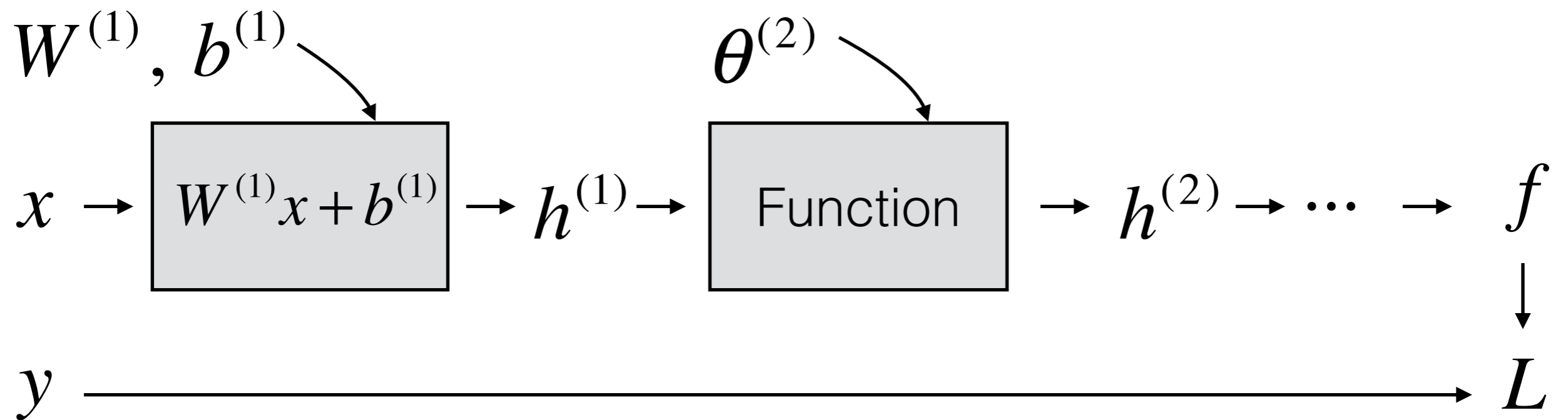


Toy Example:

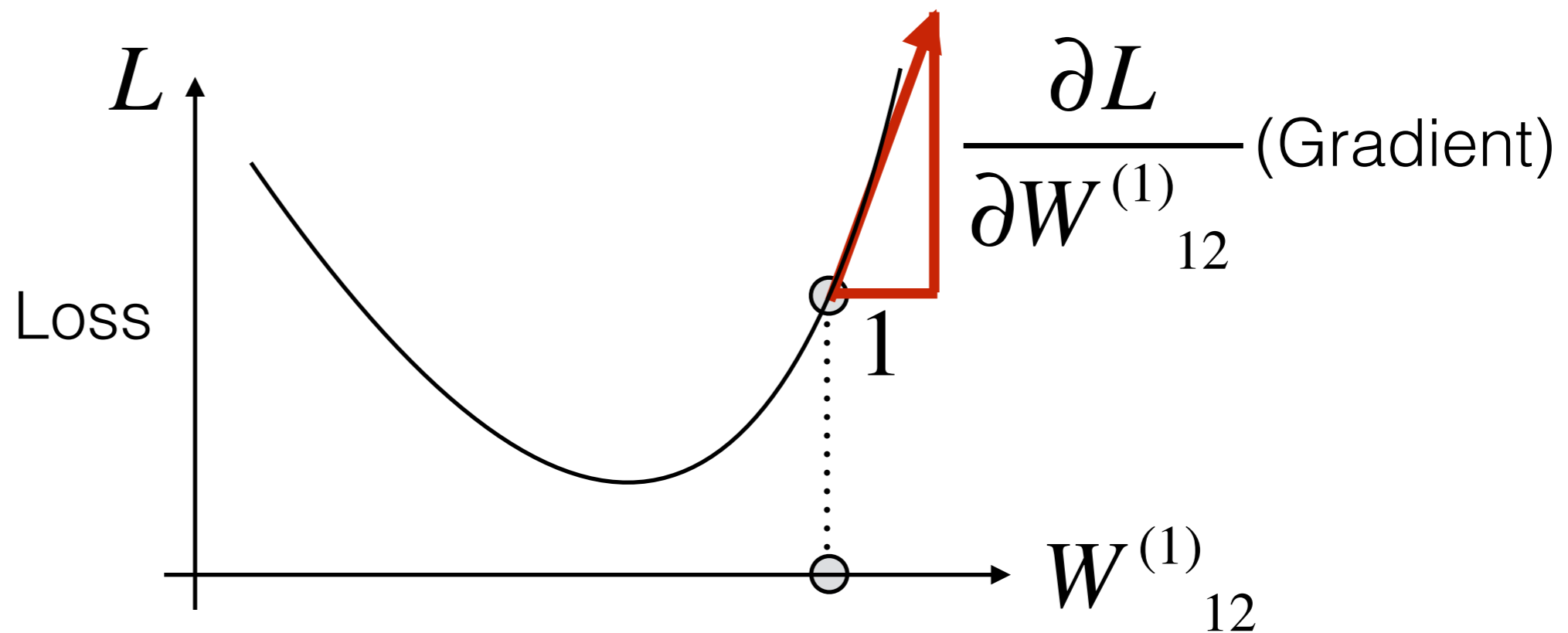


A weight somewhere in the network

Review: Setup

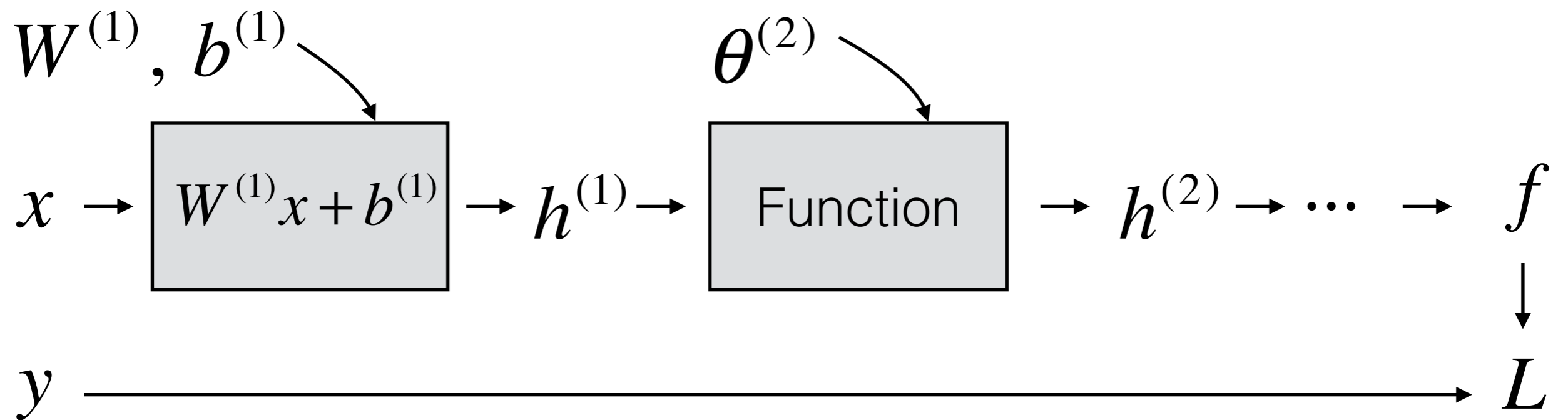


Toy Example:

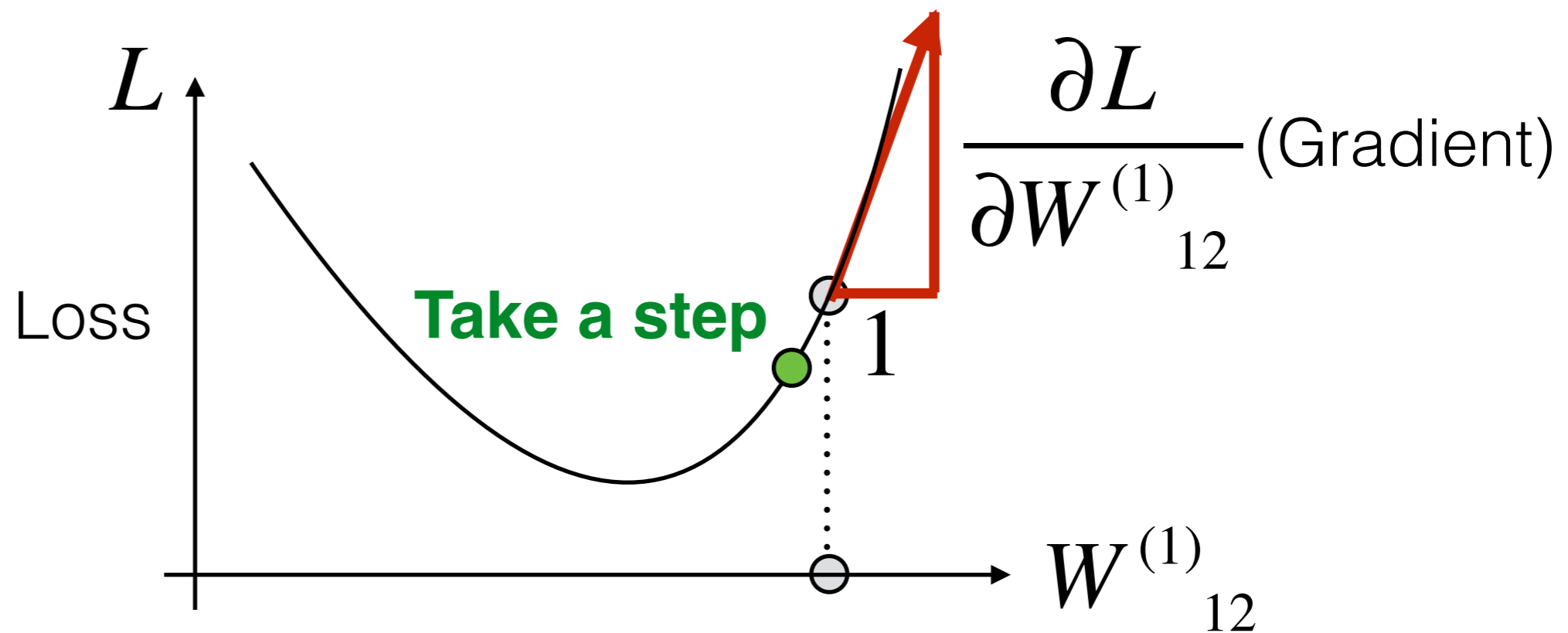


A weight somewhere in the network

Review: Setup

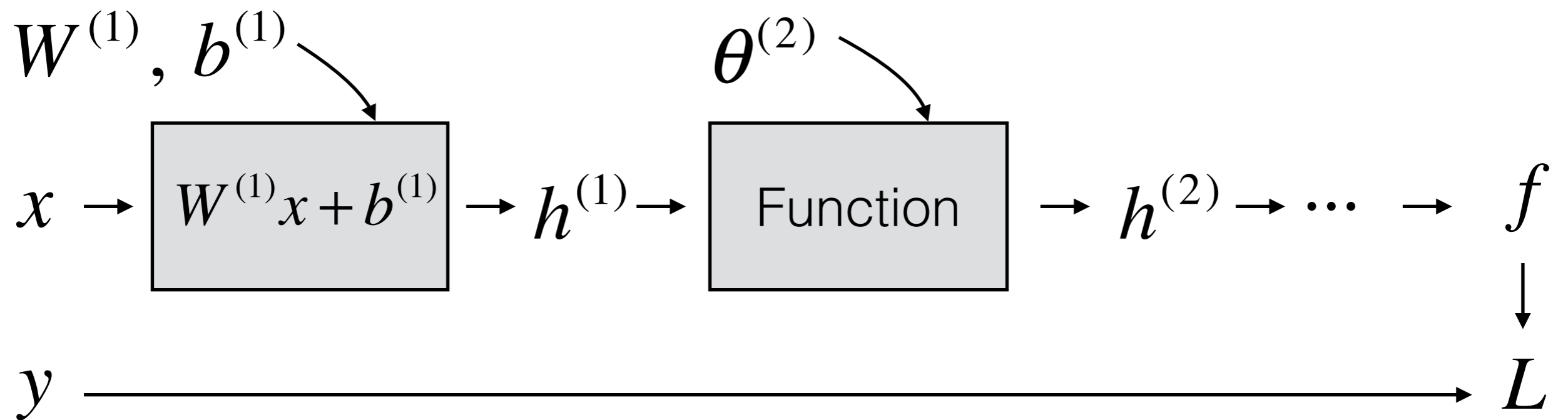


Toy Example:



A weight somewhere in the network

Review: Setup



Toy Example:

L

$\frac{\partial L}{\partial W^{(1)}}$ (Gradient)

How do we get the gradient? **Backpropagation**

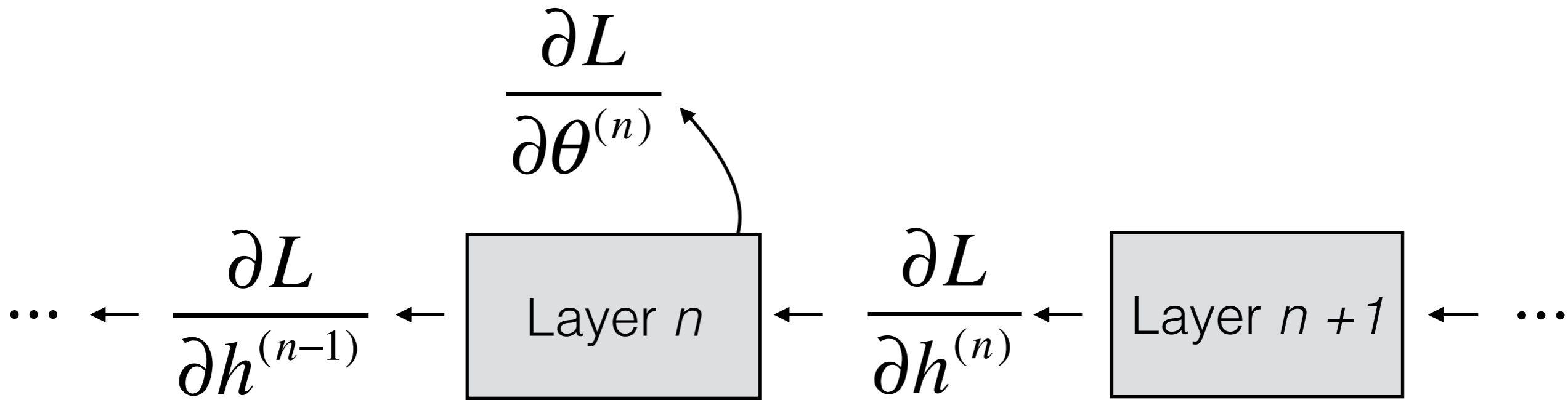
$W^{(1)}_{12}$

A weight somewhere in the network

Backprop

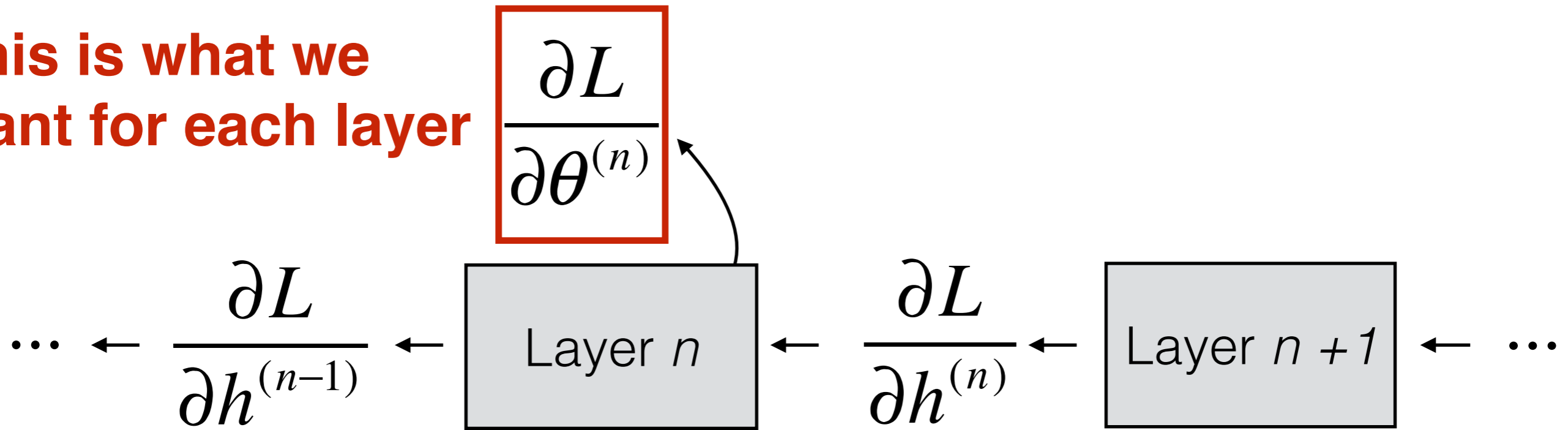
It's just the chain rule

Backprop



Backprop

This is what we want for each layer

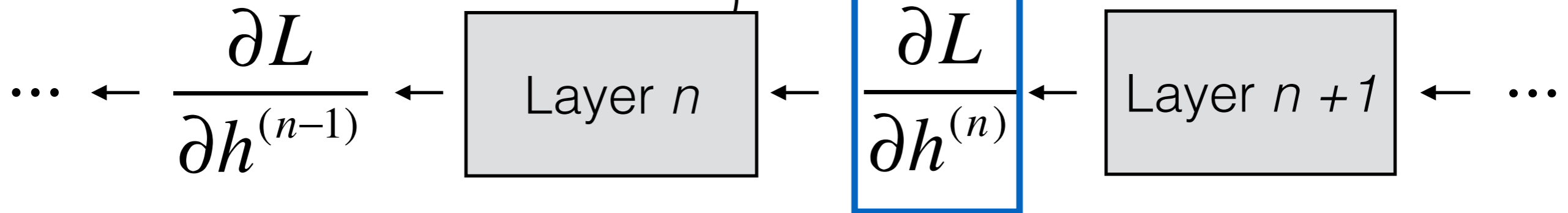


Backprop

This is what we want for each layer

$$\frac{\partial L}{\partial \theta^{(n)}}$$

To compute it, we need to propagate this gradient

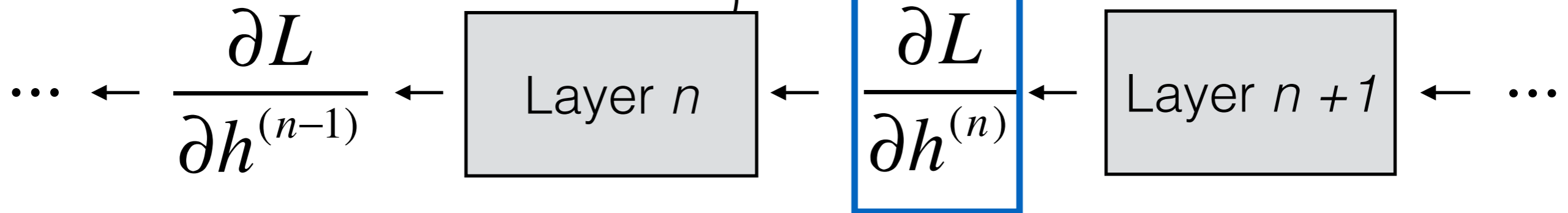


Backprop

This is what we want for each layer

$$\frac{\partial L}{\partial \theta^{(n)}}$$

To compute it, we need to propagate this gradient



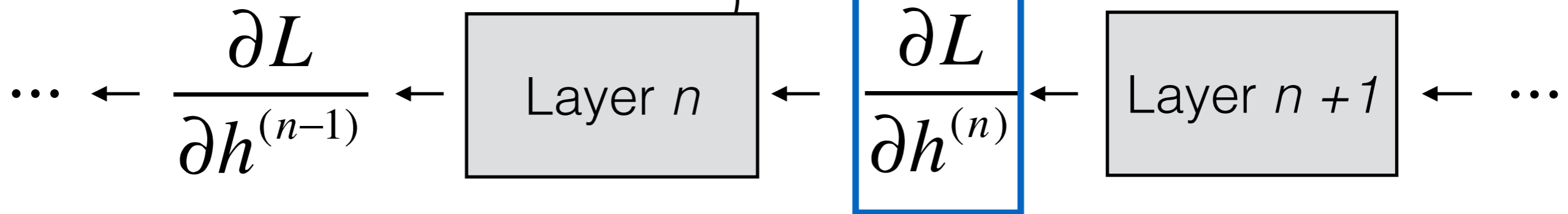
For each layer:

Backprop

This is what we want for each layer

$$\frac{\partial L}{\partial \theta^{(n)}}$$

To compute it, we need to propagate this gradient



For each layer:

$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}}$$

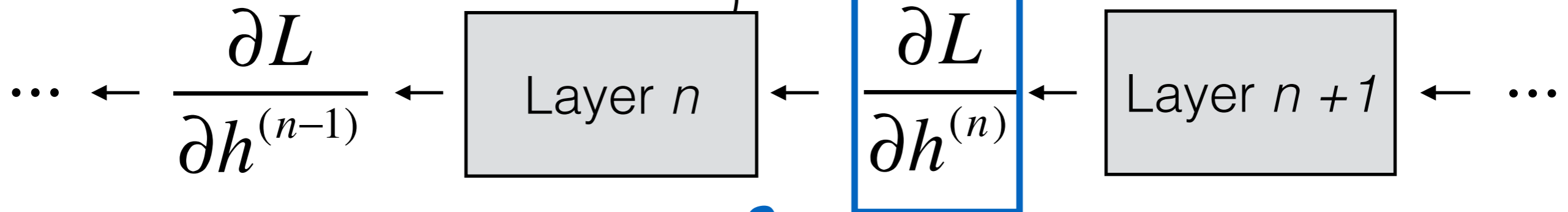
What we want

Backprop

This is what we want for each layer

$$\frac{\partial L}{\partial \theta^{(n)}}$$

To compute it, we need to propagate this gradient



For each layer:

given to us

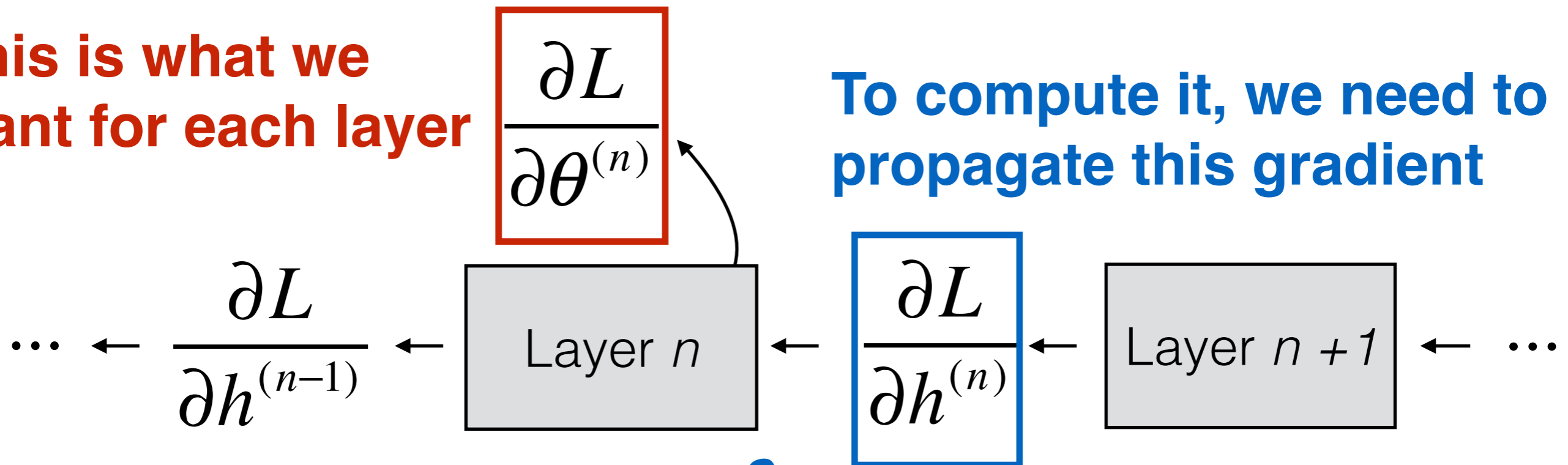
$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}}$$

What we want

Backprop

This is what we want for each layer

To compute it, we need to propagate this gradient



For each layer:

given to us

$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}}$$

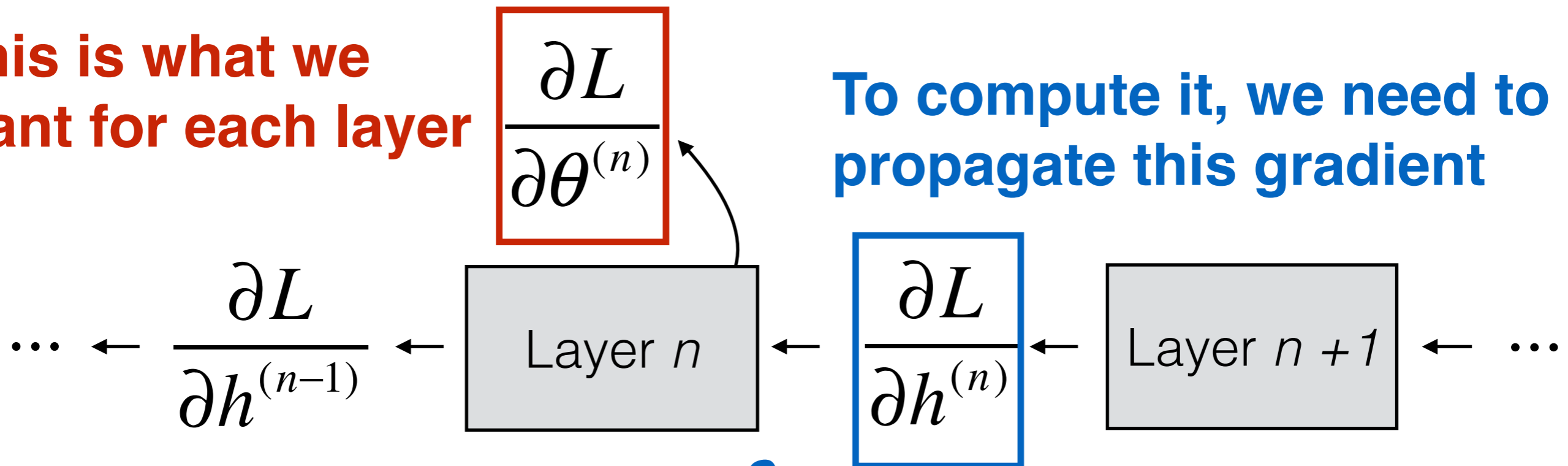
What we want

This is just the local gradient of layer n

Backprop

This is what we want for each layer

To compute it, we need to propagate this gradient



For each layer:

given to us

$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}}$$

$$\frac{\partial L}{\partial h^{(n-1)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial h^{(n-1)}}$$

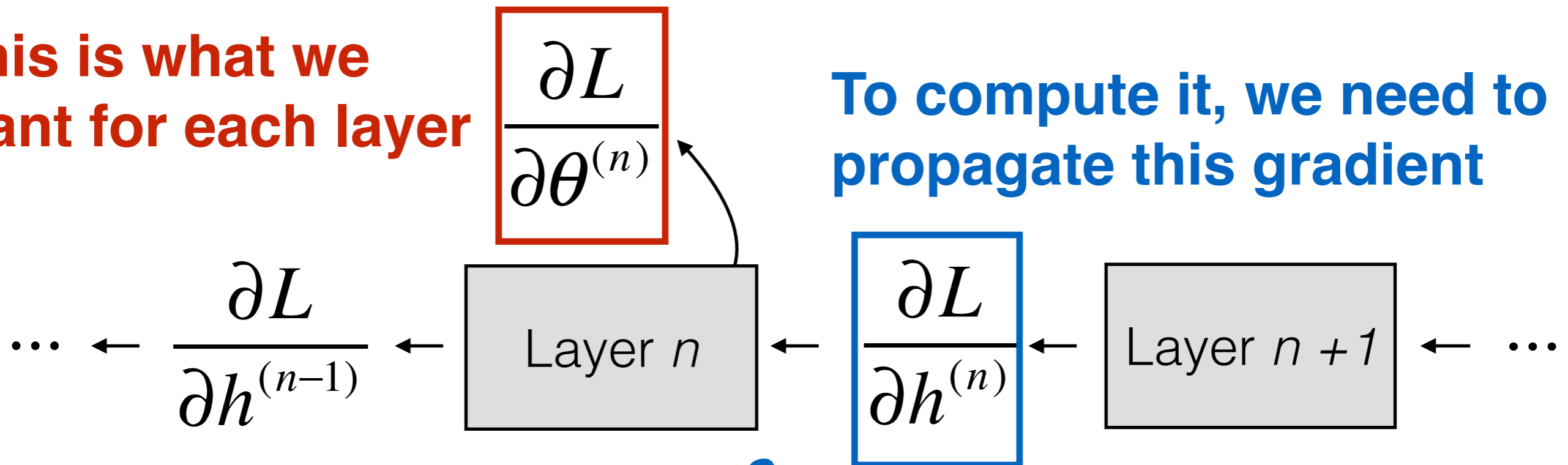
What we want

This is just the local gradient of layer n

Backprop

This is what we want for each layer

To compute it, we need to propagate this gradient



For each layer:

given to us

$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}}$$

$$\frac{\partial L}{\partial h^{(n-1)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial h^{(n-1)}}$$

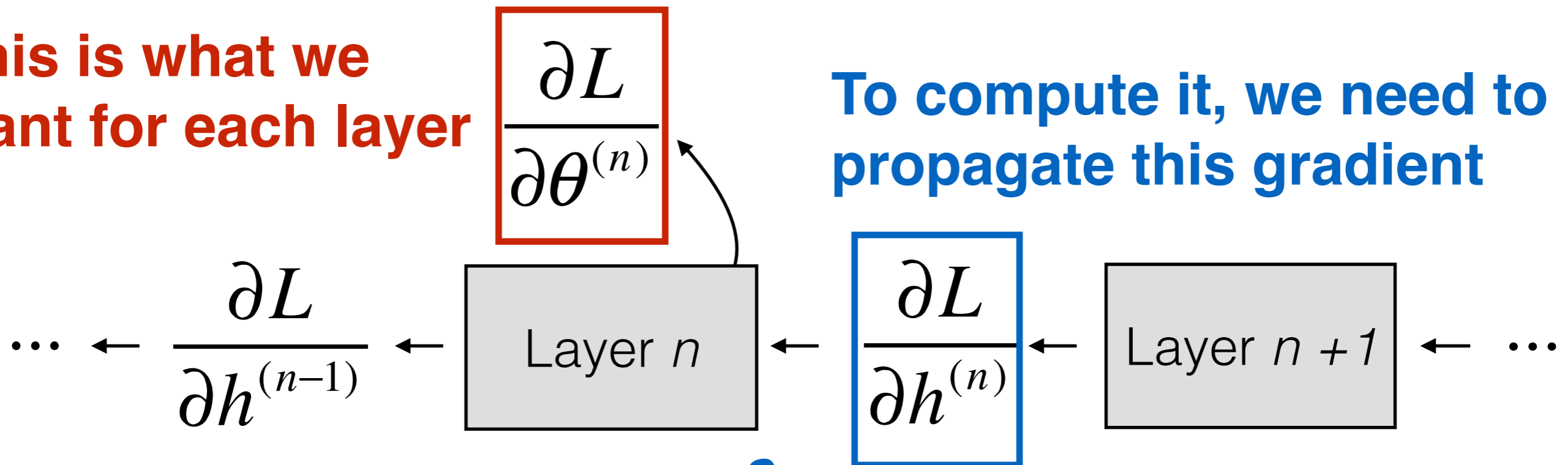
What we want

This is just the local gradient of layer n

Backprop

This is what we want for each layer

To compute it, we need to propagate this gradient



For each layer:

given to us

$$\frac{\partial L}{\partial \theta^{(n)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial \theta^{(n)}} \qquad \frac{\partial L}{\partial h^{(n-1)}} = \frac{\partial L}{\partial h^{(n)}} \cdot \frac{\partial h^{(n)}}{\partial h^{(n-1)}}$$

What we want

This is just the local gradient of layer n

Backprop

For each layer, we compute:


$$\begin{aligned} [\text{Propagated gradient to the left}] = \\ [\text{Propagated gradient from right}] \cdot [\text{Local gradient}] \end{aligned}$$

Backprop

For each layer, we compute:

$$[\text{Propagated gradient to the left}] = [\text{Propagated gradient from right}] \cdot [\text{Local gradient}]$$

(Can compute immediately)



Backprop

For each layer, we compute:

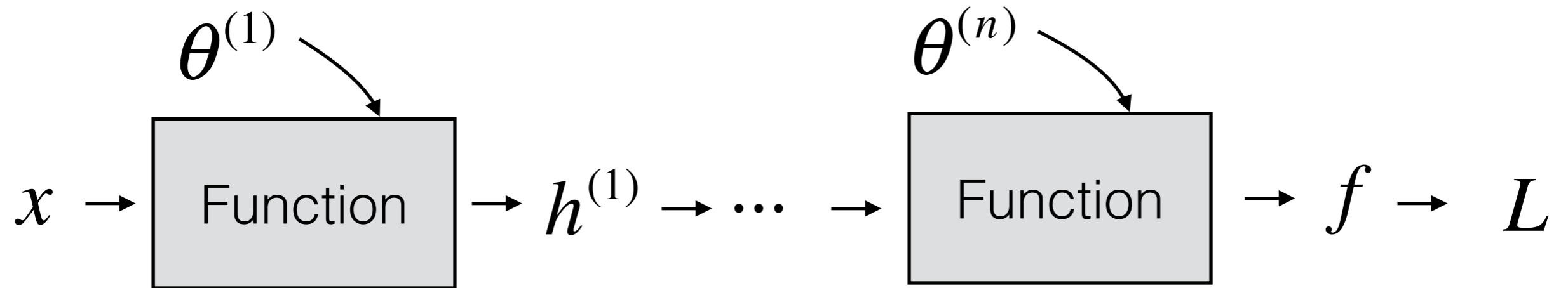
$$[\text{Propagated gradient to the left}] = [\text{Propagated gradient from right}] \cdot [\text{Local gradient}]$$

(Received during backprop) (Can compute immediately)



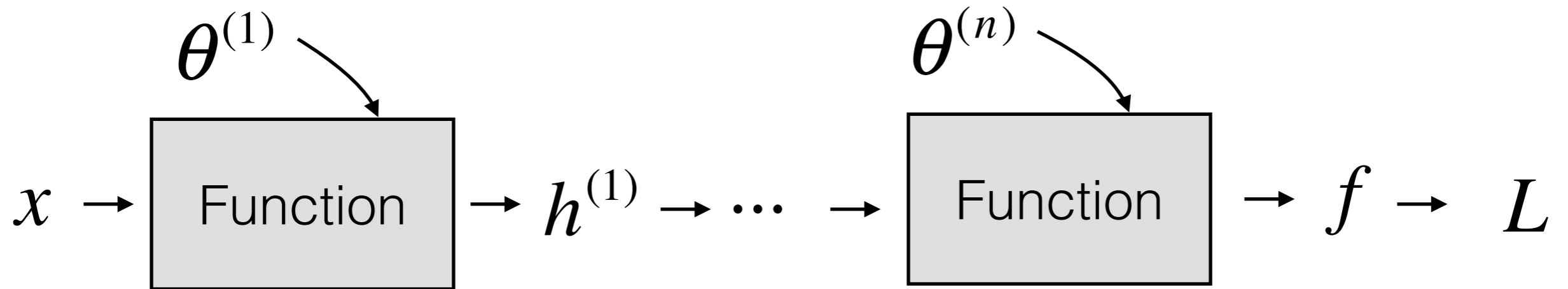
Backprop

Forward Propagation:



Backprop

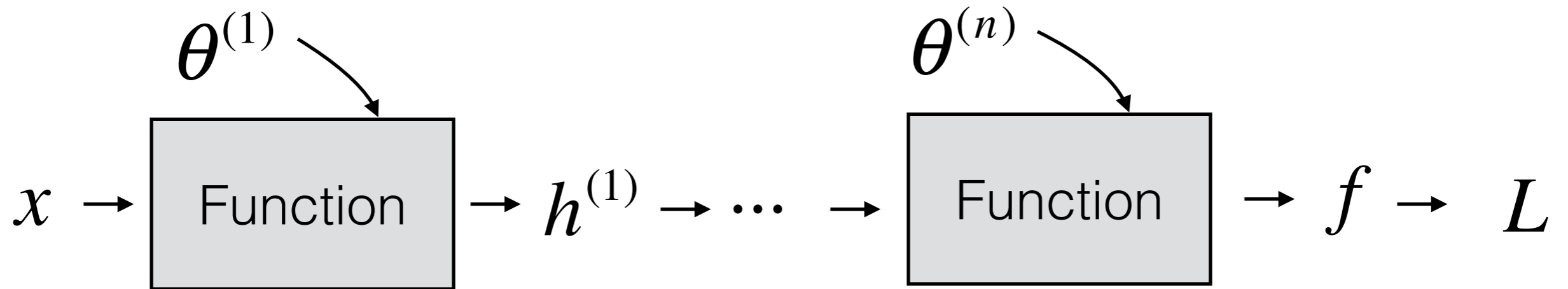
Forward Propagation:



Backward Propagation:

Backprop

Forward Propagation:

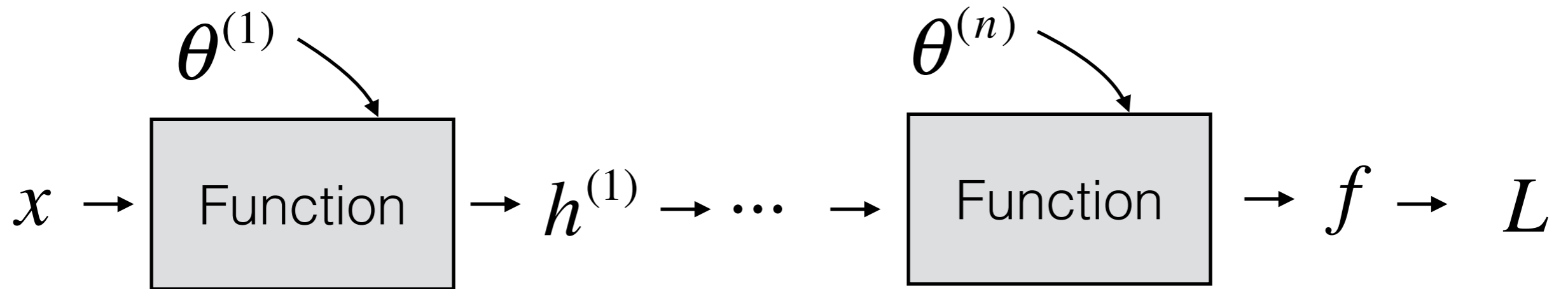


Backward Propagation:

L

Backprop

Forward Propagation:

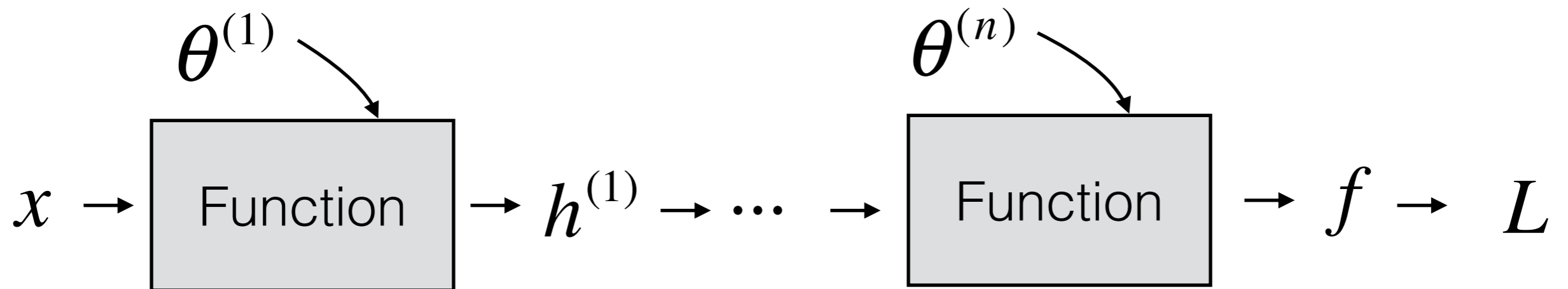


Backward Propagation:

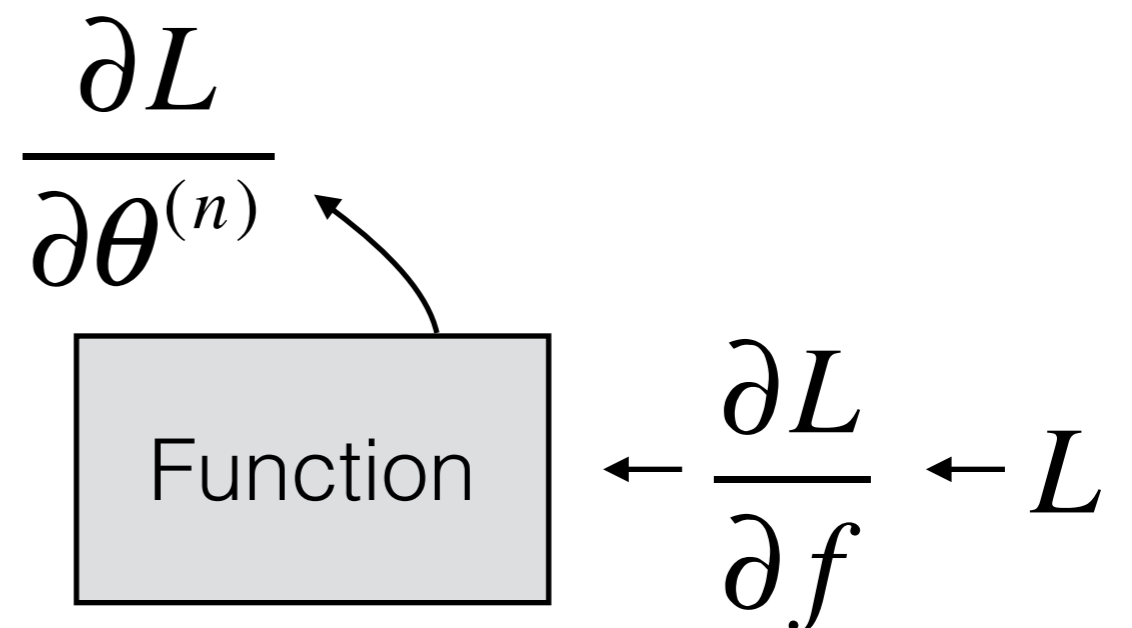
$$\frac{\partial L}{\partial f} \leftarrow L$$

Backprop

Forward Propagation:

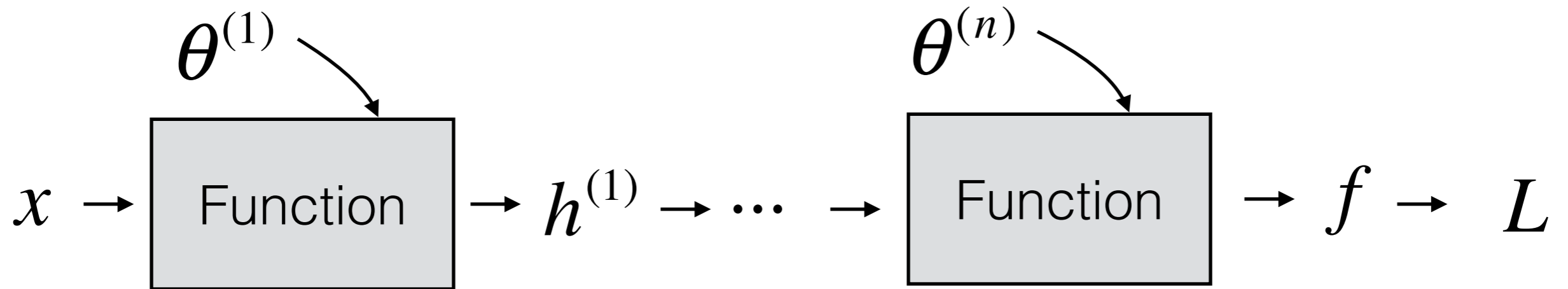


Backward Propagation:

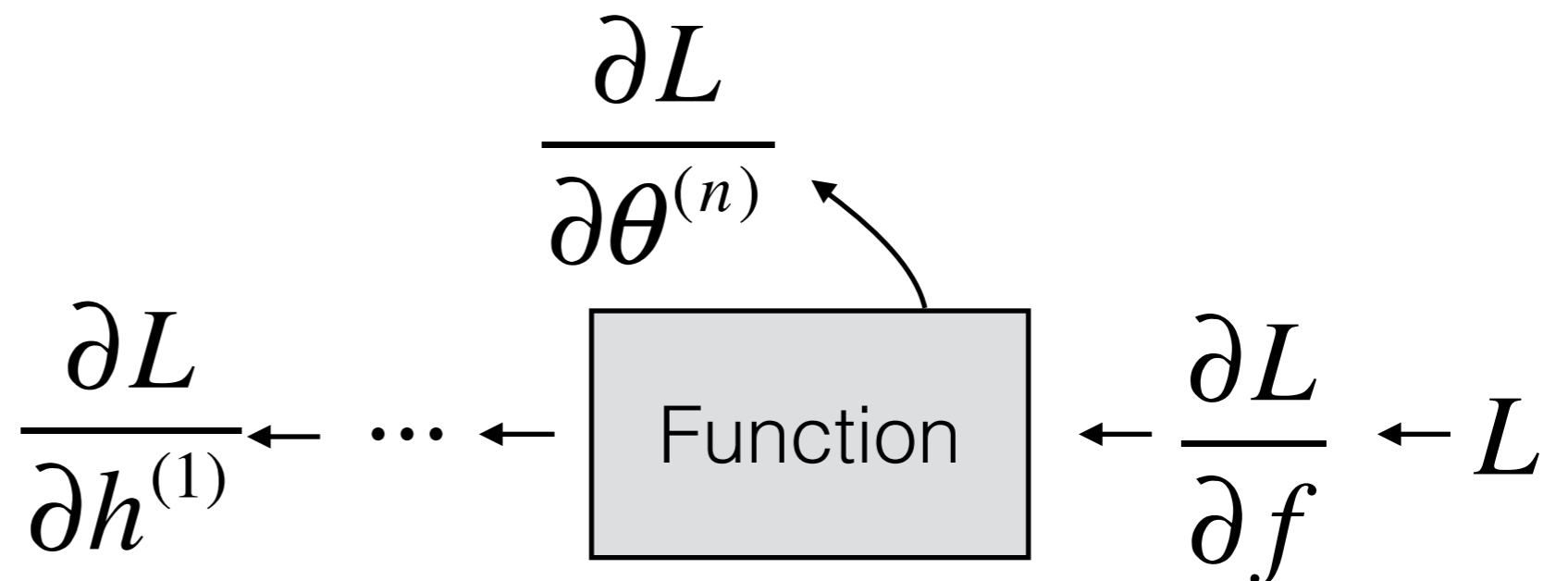


Backprop

Forward Propagation:

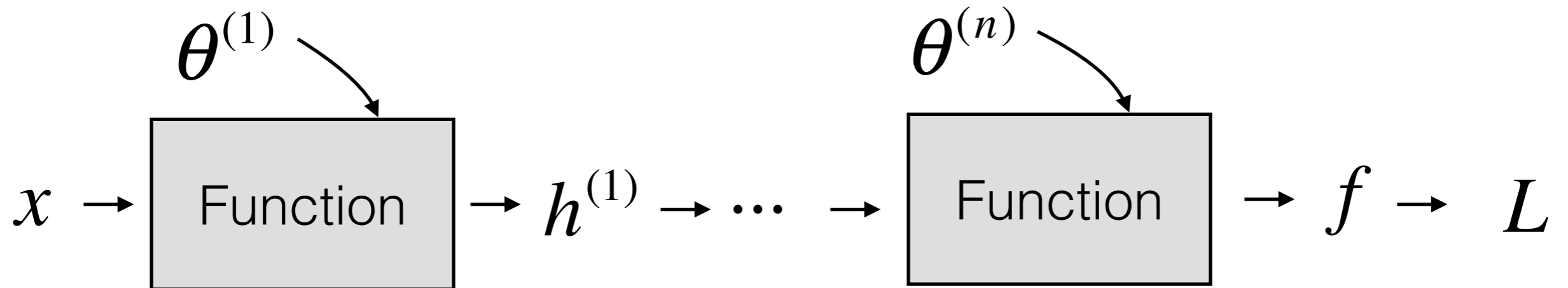


Backward Propagation:

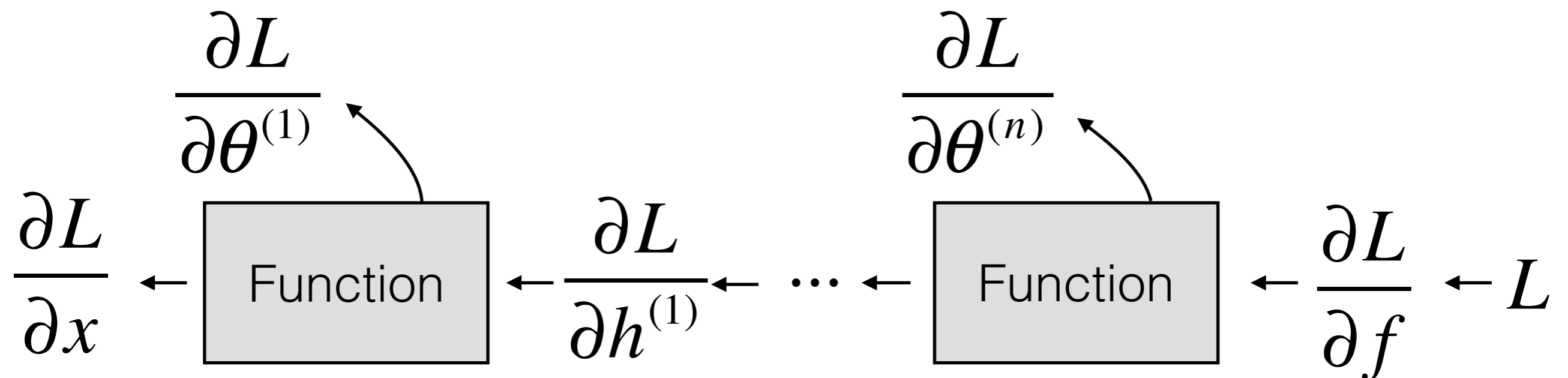


Backprop

Forward Propagation:



Backward Propagation:



Backprop

It's easy to write down the chain rule for higher dimensions —
just add more subscripts and more summations

Backprop

It's easy to write down the chain rule for higher dimensions —
just add more subscripts and more summations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial x}$$

x, h scalars
(L is always scalar)

Backprop

It's easy to write down the chain rule for higher dimensions —
just add more subscripts and more summations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial x}$$

x, h scalars
(L is always scalar)

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j}$$

x, h 1D arrays (vectors)

Backprop

It's easy to write down the chain rule for higher dimensions —
just add more subscripts and more summations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial x}$$

x, h scalars
(L is always scalar)

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j}$$

x, h 1D arrays (vectors)

$$\frac{\partial L}{\partial x_{ab}} = \sum_i \sum_j \frac{\partial L}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial x_{ab}}$$

x, h 2D arrays

Backprop

It's easy to write down the chain rule for higher dimensions —
just add more subscripts and more summations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial h} \frac{\partial h}{\partial x}$$

x, h scalars
(L is always scalar)

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j}$$

x, h 1D arrays (vectors)

$$\frac{\partial L}{\partial x_{ab}} = \sum_i \sum_j \frac{\partial L}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial x_{ab}}$$

x, h 2D arrays

$$\frac{\partial L}{\partial x_{abc}} = \sum_i \sum_j \sum_k \frac{\partial L}{\partial h_{ijk}} \frac{\partial h_{ijk}}{\partial x_{abc}}$$

x, h 3D arrays

Example: Mean Subtraction (for a single input)

Example: Mean Subtraction

(for a single input)

- Ok, so how do we actually derive the backwards pass? Let's walk through an example together.

Example: Mean Subtraction

(for a single input)

- Ok, so how do we actually derive the backwards pass? Let's walk through an example together.
- Example layer: mean subtraction:

$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Example: Mean Subtraction

(for a single input)

- Ok, so how do we actually derive the backwards pass? Let's walk through an example together.
- Example layer: mean subtraction:

$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

(here, "i" and "k"
are channels)

Example: Mean Subtraction

(for a single input)

- Ok, so how do we actually derive the backwards pass? Let's walk through an example together.

- Example layer: mean subtraction:

$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

(here, "i" and "k" are channels)

- For backprop, we just need the local derivative

Example: Mean Subtraction (for a single input)

Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$
- Taking the derivative of the layer:

Example: Mean Subtraction

(for a single input)


- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$
- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$


$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$


Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j} \quad (\text{backprop aka chain rule})$$

$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$


Example: Mean Subtraction


(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j} \quad (\text{backprop aka chain rule})$$

$$= \sum_i \frac{\partial L}{\partial h_i} \left(\delta_{ij} - \frac{1}{D} \right)$$

$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$


Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j} \quad (\text{backprop aka chain rule})$$

$$= \sum_i \frac{\partial L}{\partial h_i} \left(\delta_{ij} - \frac{1}{D} \right)$$

$$= \sum_i \frac{\partial L}{\partial h_i} \delta_{ij} - \frac{1}{D} \sum_i \frac{\partial L}{\partial h_i}$$

$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j} \quad (\text{backprop aka chain rule})$$

$$= \sum_i \frac{\partial L}{\partial h_i} \left(\delta_{ij} - \frac{1}{D} \right)$$

$$= \sum_i \frac{\partial L}{\partial h_i} \delta_{ij} - \frac{1}{D} \sum_i \frac{\partial L}{\partial h_i}$$

$$= \frac{\partial L}{\partial h_j} - \frac{1}{D} \sum_i \frac{\partial L}{\partial h_i}$$

$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

Example: Mean Subtraction

(for a single input)

- Forward: $h_i = x_i - \frac{1}{D} \sum_k x_k$

- Taking the derivative of the layer: $\frac{\partial h_i}{\partial x_j} = \delta_{ij} - \frac{1}{D}$

$$\frac{\partial L}{\partial x_j} = \sum_i \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial x_j} \quad (\text{backprop aka chain rule})$$

$$= \sum_i \frac{\partial L}{\partial h_i} \left(\delta_{ij} - \frac{1}{D} \right)$$

$$= \sum_i \frac{\partial L}{\partial h_i} \delta_{ij} - \frac{1}{D} \sum_i \frac{\partial L}{\partial h_i}$$

$$= \frac{\partial L}{\partial h_j} - \frac{1}{D} \sum_i \frac{\partial L}{\partial h_i}$$

Done!

$$\left(\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \right)$$

Example: Mean Subtraction

(for a single input)

$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial h_i} - \frac{1}{D} \sum_k \frac{\partial L}{\partial h_k}$$

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$
- Backward:
$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial h_i} - \frac{1}{D} \sum_k \frac{\partial L}{\partial h_k}$$

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$
- Backward:
$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial h_i} - \frac{1}{D} \sum_k \frac{\partial L}{\partial h_k}$$
- In this case, they're identical operations!

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$
- Backward:
$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial h_i} - \frac{1}{D} \sum_k \frac{\partial L}{\partial h_k}$$
- In this case, they're identical operations!
- Usually the forwards pass and backwards pass are similar **but not the same**.

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$
- Backward:
$$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial h_i} - \frac{1}{D} \sum_k \frac{\partial L}{\partial h_k}$$
- In this case, they're identical operations!
- Usually the forwards pass and backwards pass are similar **but not the same**.
- Derive it by hand, and check it numerically

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Let's code this up in NumPy:

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Let's code this up in NumPy:

```
def forward(X):  
    return X - np.mean(X, axis=1)
```

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Let's code this up in NumPy:

```
def forward(X):  
    return X - np.mean(X, axis=1)
```

Dimension mismatch

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Let's code this up in NumPy:

```
def forward(X):  
    return X - np.mean(X, axis=1)
```

Dimension mismatch

You need to broadcast properly:

```
def forward(X):  
    return X - np.mean(X, axis=1)[:, np.newaxis]
```

Example: Mean Subtraction

(for a single input)

- Forward:
$$h_i = x_i - \frac{1}{D} \sum_k x_k$$

Let's code this up in NumPy:

```
def forward(X): Dimension mismatch  
    return X - np.mean(X, axis=1)
```

You need to broadcast properly:

```
def forward(X):  
    return X - np.mean(X, axis=1)[:, np.newaxis]
```

This also works:

```
def forward(X):  
    return X - np.mean(X, axis=1, keepdims=True)
```

Example: Mean Subtraction

(for a single input)

The backward pass is easy:

```
def backward(dh):  
    return forward(dh)
```

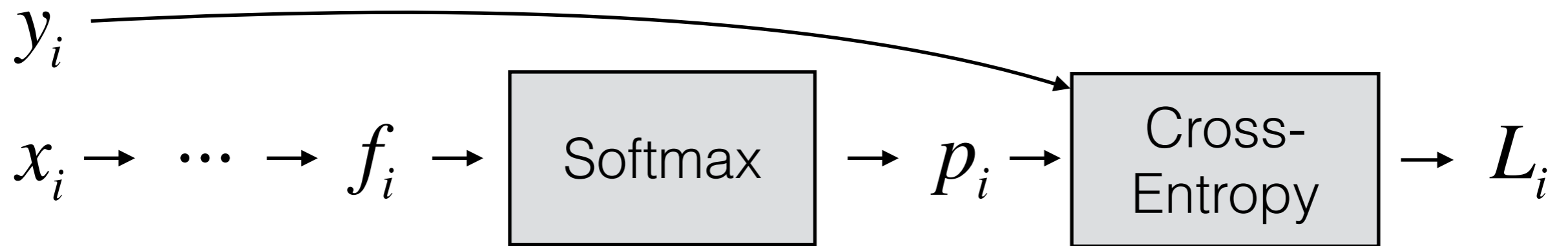
(Remember they're usually not the same)

Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss
(together this is often called "Softmax loss")

Example: Softmax (for N inputs)

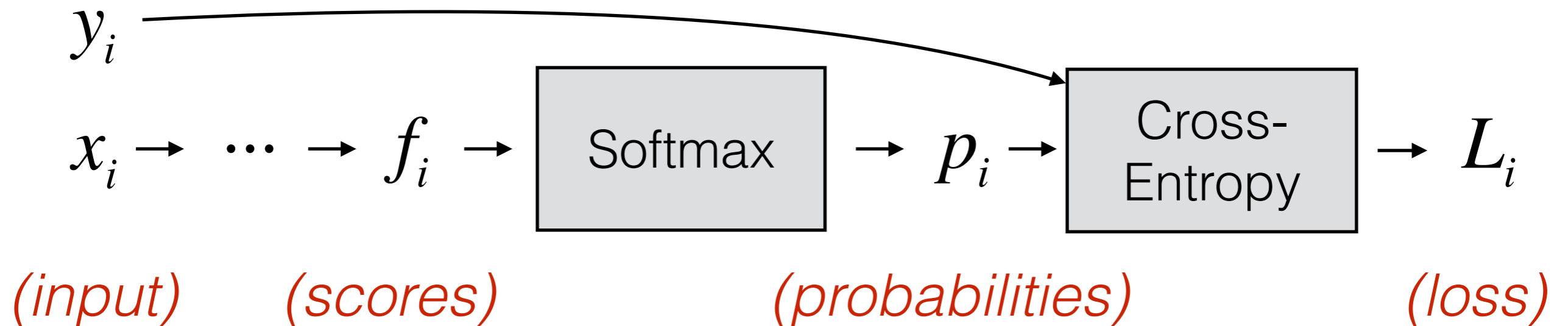
Let's assume we are using Softmax and Cross-entropy loss (together this is often called "Softmax loss")



Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss
(together this is often called "Softmax loss")

(ground truth labels)

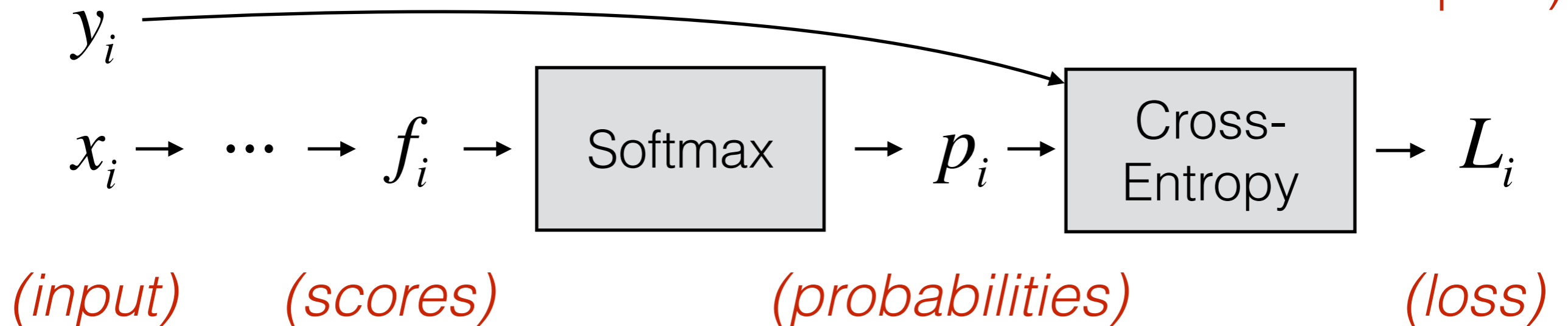


Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss (together this is often called "Softmax loss")

(ground truth labels)

(here, "i" are different examples)

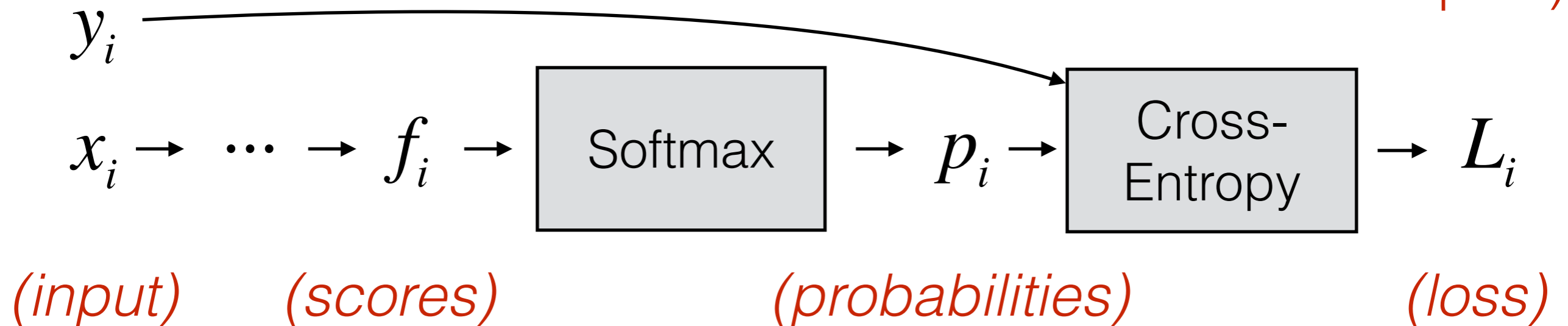


Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss (together this is often called "Softmax loss")

(ground truth labels)

(here, "i" are different examples)



$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

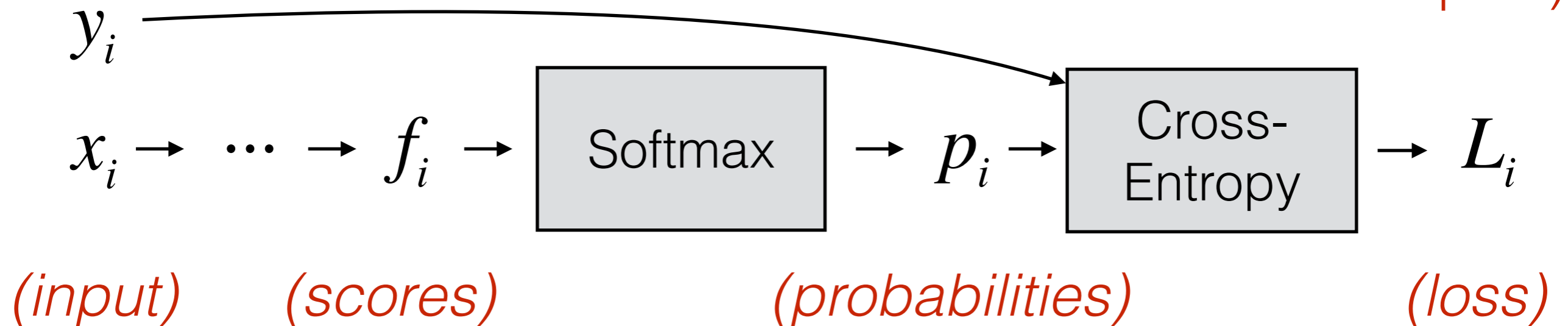
(Softmax)

Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss (together this is often called "Softmax loss")

(ground truth labels)

(here, "i" are different examples)



$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

(Softmax)

$$L_i = -\log p_{i,y_i}$$

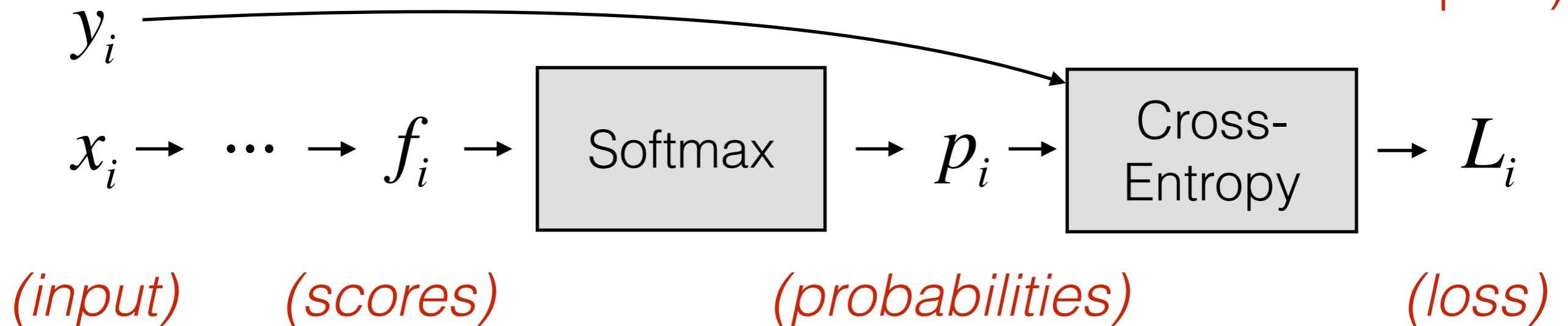
(Cross-entropy)

Example: Softmax (for N inputs)

Let's assume we are using Softmax and Cross-entropy loss (together this is often called "Softmax loss")

(ground truth labels)

(here, "i" are different examples)



$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

(Softmax)

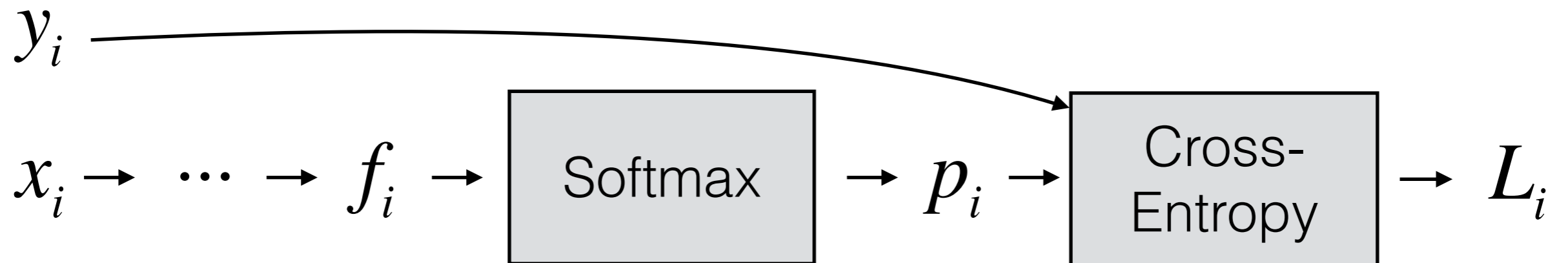
$$L_i = -\log p_{i,y_i}$$

(Cross-entropy)

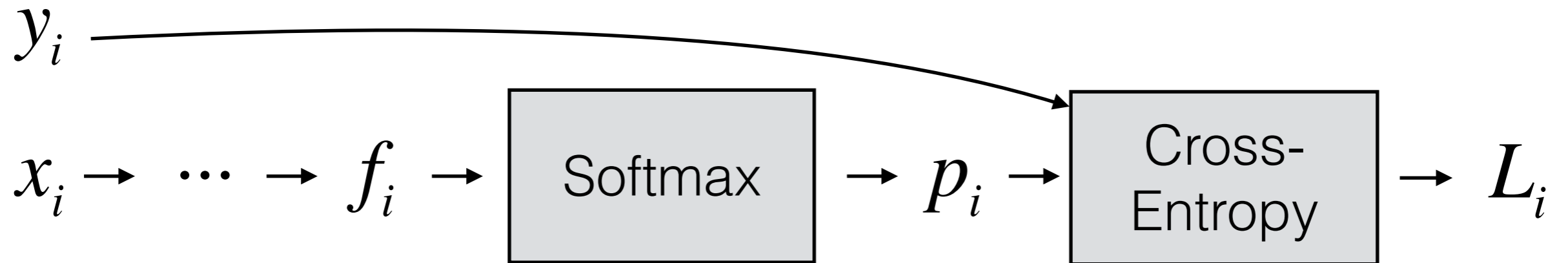
$$L = \frac{1}{N} \sum_i L_i$$

(Avg. over examples)

Example: Softmax (for N inputs)

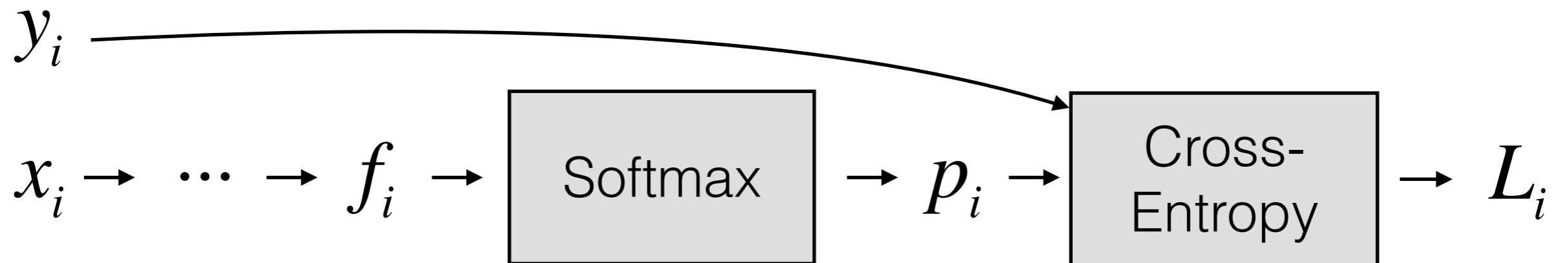


Example: Softmax (for N inputs)



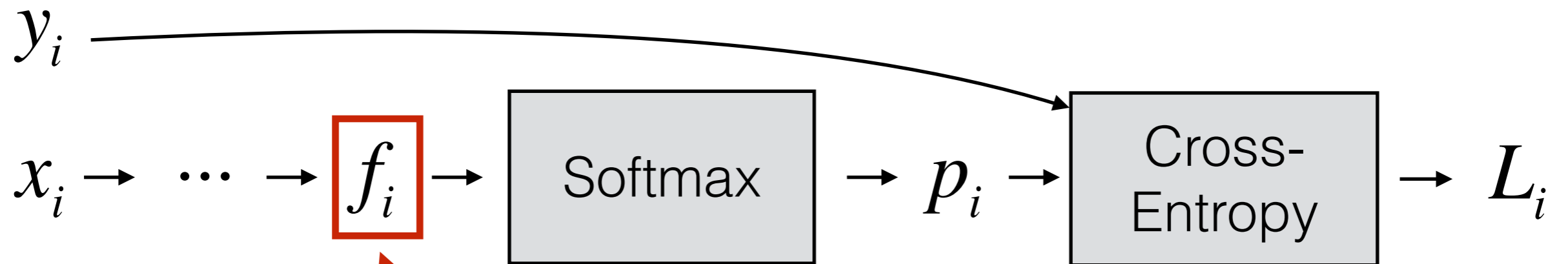
Derivative:
$$\frac{\partial L}{\partial f_{i,j}} = \frac{p_{i,j} - t_{i,j}}{N}$$

Example: Softmax (for N inputs)



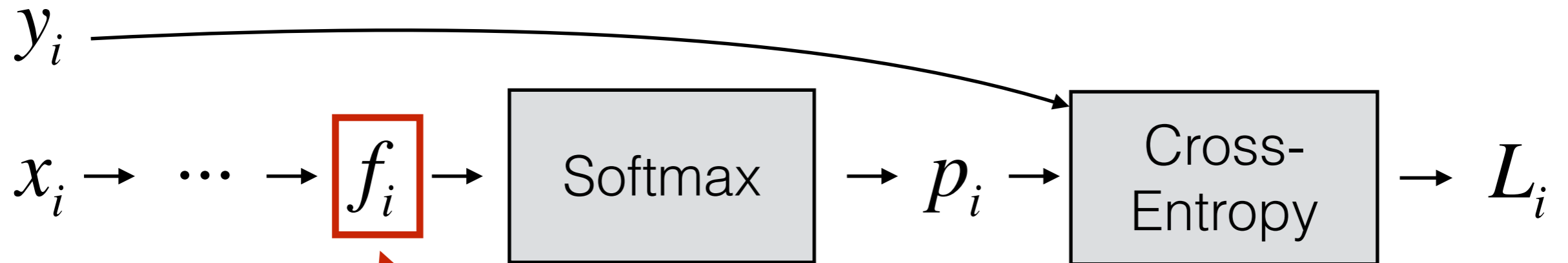
Derivative: $\frac{\partial L}{\partial f_{i,j}} = \frac{p_{i,j} - t_{i,j}}{N}$ where $t_i = [0 \dots 1 \dots 0]$
(Entry y_i set to 1)

Example: Softmax (for N inputs)



Derivative: $\frac{\partial L}{\partial f_{i,j}} = \frac{p_{i,j} - t_{i,j}}{N}$ where $t_i = [0 \dots 1 \dots 0]$
(Entry y_i set to 1)

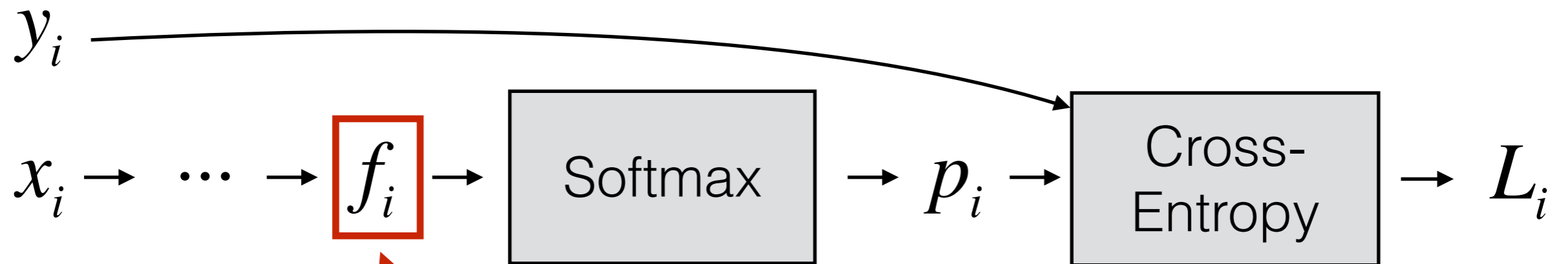
Example: Softmax (for N inputs)



Derivative: $\frac{\partial L}{\partial f_{i,j}} = \frac{p_{i,j} - t_{i,j}}{N}$ where $t_i = [0 \dots 1 \dots 0]$
(Entry y_i set to 1)

(Try deriving this — it's tricky but not too hard)

Example: Softmax (for N inputs)



Derivative: $\frac{\partial L}{\partial f_{i,j}} = \frac{p_{i,j} - t_{i,j}}{N}$ where $t_i = [0 \dots 1 \dots 0]$
(Entry y_i set to 1)

(Try deriving this — it's tricky but not too hard)

Now we can continue backpropagating to the layer before “f”

Example: Softmax (for N inputs)

Let's code this up in NumPy:

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

Example: Softmax (for N inputs)

Let's code this up in NumPy:

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

Doesn't work — what's the problem this time?

Example: Softmax (for N inputs)

Let's code this up in NumPy:

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

Doesn't work — what's the problem this time?

- What if there is the value 1000 appears in "f"?

Example: Softmax (for N inputs)

Let's code this up in NumPy:

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

Doesn't work — what's the problem this time?

- What if there is the value 1000 appears in "f"?

Overflow —> *we get inf/inf = NaN*

Example: Softmax (for N inputs)

Let's code this up in NumPy:

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

Doesn't work — what's the problem this time?

- What if there is the value 1000 appears in "f"?

Overflow —> *we get inf/inf = NaN*

- What if the largest value is -1000?

Example: Softmax (for N inputs)

Let's code this up in NumPy:

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

Doesn't work — what's the problem this time?

- What if there is the value 1000 appears in "f"?

Overflow —> we get $inf/inf = NaN$

- What if the largest value is -1000?

Underflow —> we get $0/0 = NaN$

Example: Softmax (for N inputs)

Let's code this up in NumPy:

$$p_{i,j} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

```
def softmax(f):  
    exp_f = np.exp(f)  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

Doesn't work — what's the problem this time?

- What if there is the value 1000 appears in "f"?

Overflow —> we get $inf/inf = NaN$

- What if the largest value is -1000?

Underflow —> we get $0/0 = NaN$

This expression is numerically unstable

Example: Softmax (for N inputs)

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

$$p_{i,j} = \frac{e^{f_{i,j}-C}}{\sum_k e^{f_{j,k}-C}} = \frac{e^{-C} e^{f_{i,j}}}{\sum_k e^{-C} e^{f_{i,k}}} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

$$P_{i,j} = \frac{e^{f_{i,j}-C}}{\sum_k e^{f_{j,k}-C}} = \frac{e^{-C} e^{f_{i,j}}}{\sum_k e^{-C} e^{f_{i,k}}} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

If we choose “C” to be the max, then it works:

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

$$P_{i,j} = \frac{e^{f_{i,j}-C}}{\sum_k e^{f_{j,k}-C}} = \frac{e^{-C} e^{f_{i,j}}}{\sum_k e^{-C} e^{f_{i,k}}} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

If we choose “C” to be the max, then it works:

- If a large value appears in “f”, then that value will become 1 and all others will be 0 (avoiding overflow)

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

$$P_{i,j} = \frac{e^{f_{i,j}-C}}{\sum_k e^{f_{j,k}-C}} = \frac{e^{-C} e^{f_{i,j}}}{\sum_k e^{-C} e^{f_{i,k}}} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

If we choose “C” to be the max, then it works:

- If a large value appears in “f”, then that value will become 1 and all others will be 0 (avoiding overflow)
- If all values in “f” are large negative, then they will be shifted up towards 0 (avoiding underflow)

Example: Softmax (for N inputs)

Observation: subtracting a constant does not change “p”:

$$P_{i,j} = \frac{e^{f_{i,j}-C}}{\sum_k e^{f_{j,k}-C}} = \frac{e^{-C} e^{f_{i,j}}}{\sum_k e^{-C} e^{f_{i,k}}} = \frac{e^{f_{i,j}}}{\sum_k e^{f_{i,k}}}$$

If we choose “C” to be the max, then it works:

- If a large value appears in “f”, then that value will become 1 and all others will be 0 (avoiding overflow)
- If all values in “f” are large negative, then they will be shifted up towards 0 (avoiding underflow)

```
def softmax(f):  
    exp_f = np.exp(f - np.max(f, axis=1, keepdims=True))  
    return exp_f / np.sum(exp_f, axis=1, keepdims=True)
```

What about the weights?

To get the derivative of the weights, use the chain rule again!

What about the weights?

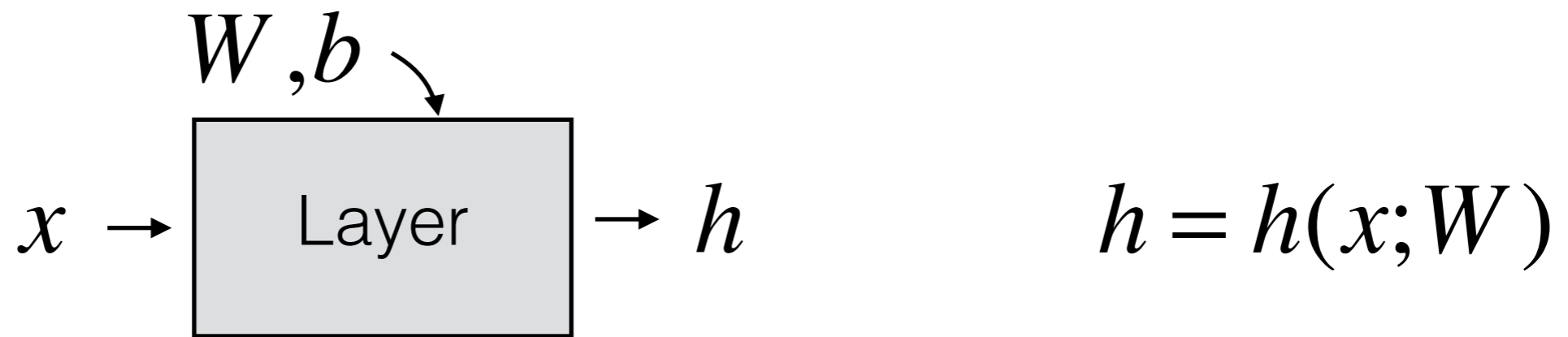
To get the derivative of the weights, use the chain rule again!

Example: 2D weights, 1D bias, 1D hidden activations:

What about the weights?

To get the derivative of the weights, use the chain rule again!

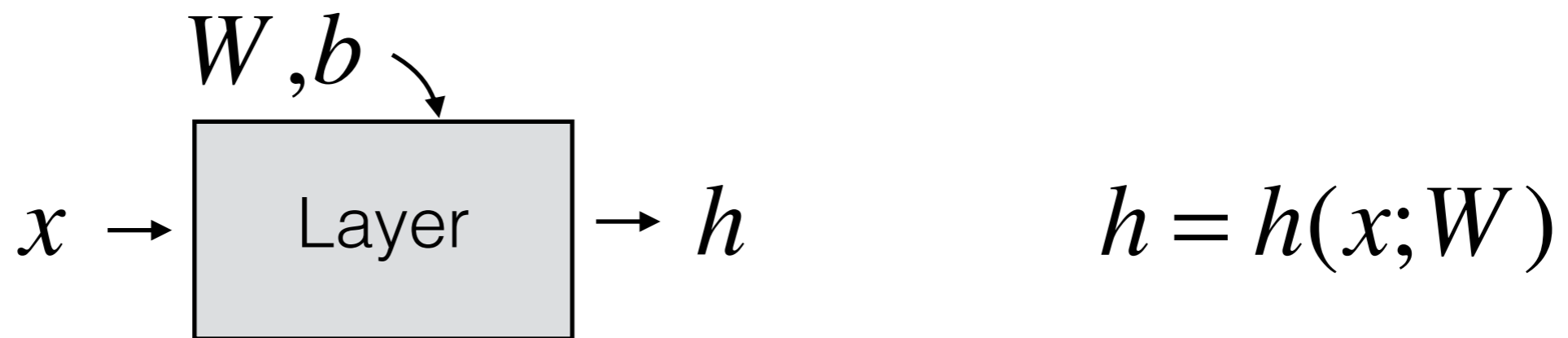
Example: 2D weights, 1D bias, 1D hidden activations:



What about the weights?

To get the derivative of the weights, use the chain rule again!

Example: 2D weights, 1D bias, 1D hidden activations:

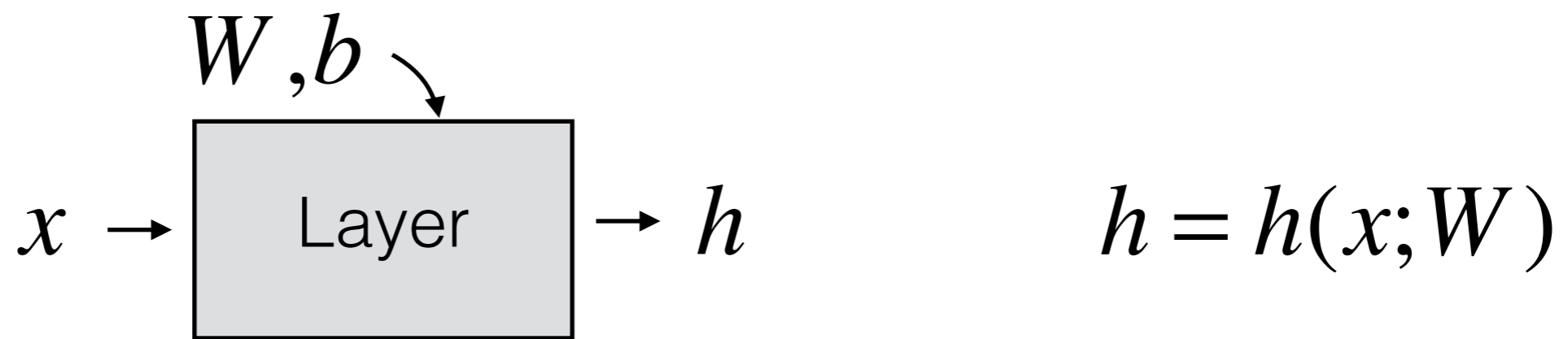


$$\frac{\partial L}{\partial W_{ij}} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial W_{ij}}$$

What about the weights?

To get the derivative of the weights, use the chain rule again!

Example: 2D weights, 1D bias, 1D hidden activations:



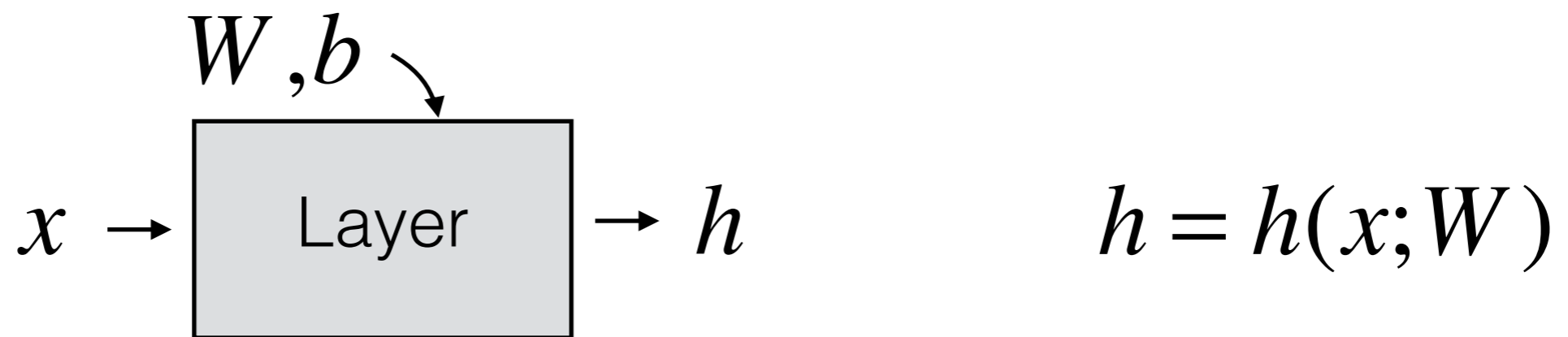
$$\frac{\partial L}{\partial W_{ij}} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial W_{ij}}$$

$$\frac{\partial L}{\partial b_i} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial b_i}$$

What about the weights?

To get the derivative of the weights, use the chain rule again!

Example: 2D weights, 1D bias, 1D hidden activations:



$$\frac{\partial L}{\partial W_{ij}} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial W_{ij}}$$

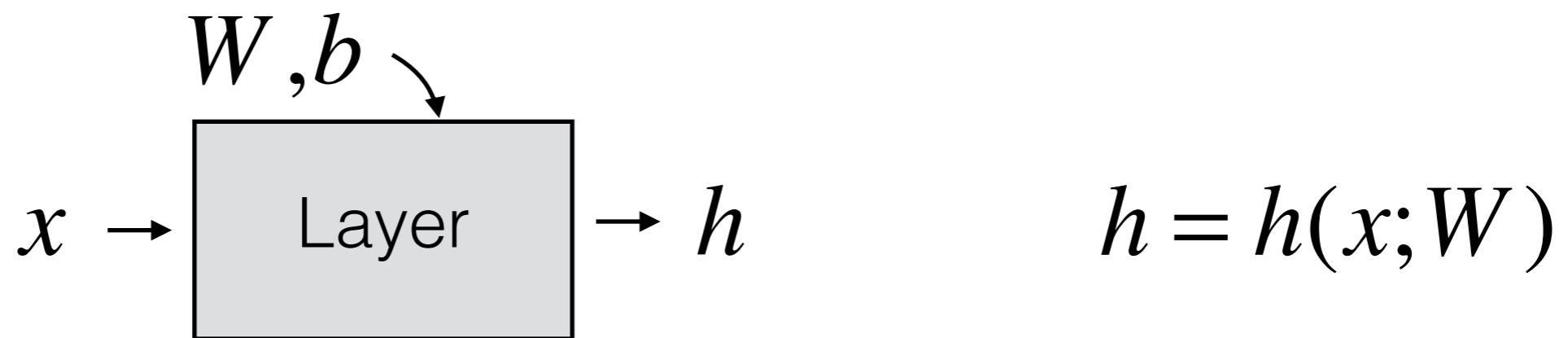
$$\frac{\partial L}{\partial b_i} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial b_i}$$

(the number of subscripts and summations changes depending on your layer and parameter sizes)

What about the weights?

To get the derivative of the weights, use the chain rule again!

Example: 2D weights, 1D bias, 1D hidden activations:



$$\frac{\partial L}{\partial W_{ij}} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial W_{ij}}$$

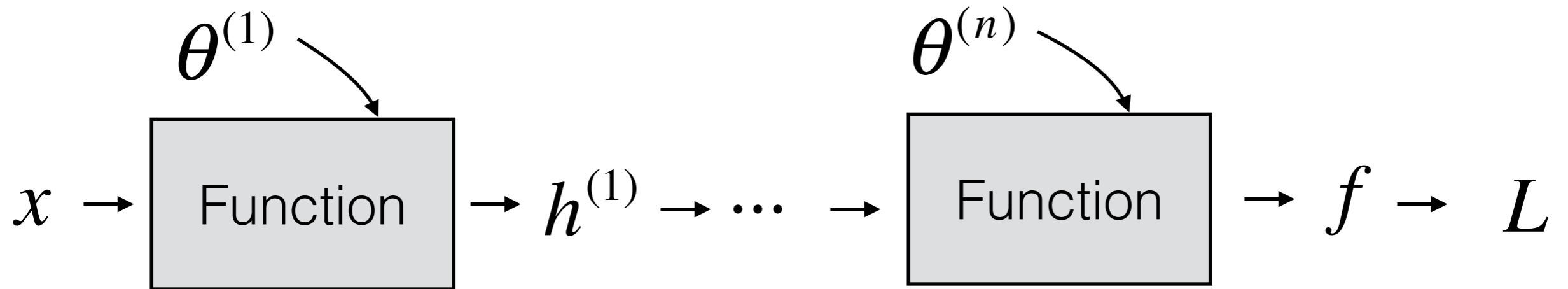
$$\frac{\partial L}{\partial b_i} = \sum_k \frac{\partial L}{\partial h_k} \frac{\partial h_k}{\partial b_i}$$

(the number of subscripts and summations changes depending on your layer and parameter sizes)

HW2: you will derive this for various layers.

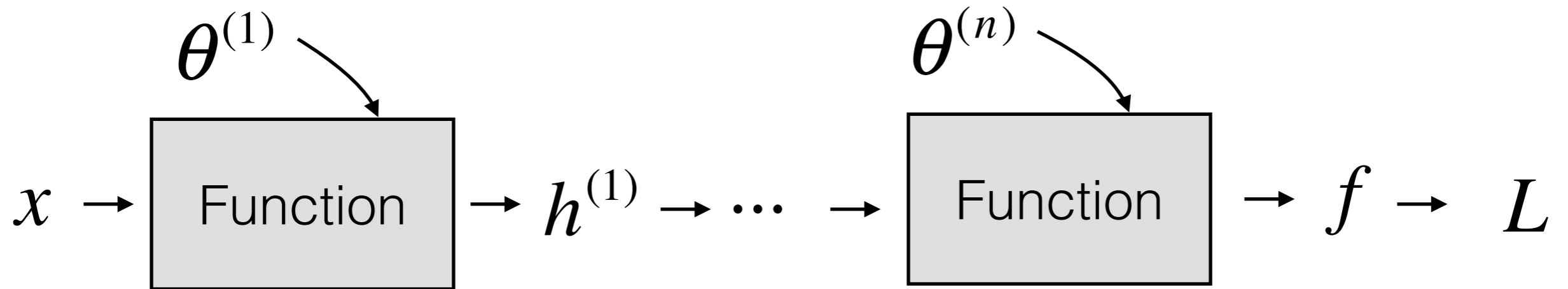
Recap

Forward Propagation:



Recap

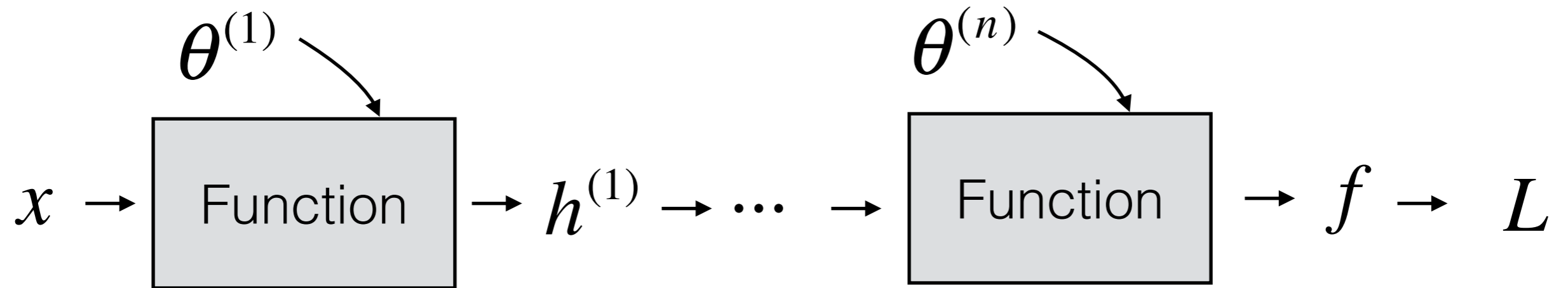
Forward Propagation:



Backward Propagation:

Recap

Forward Propagation:

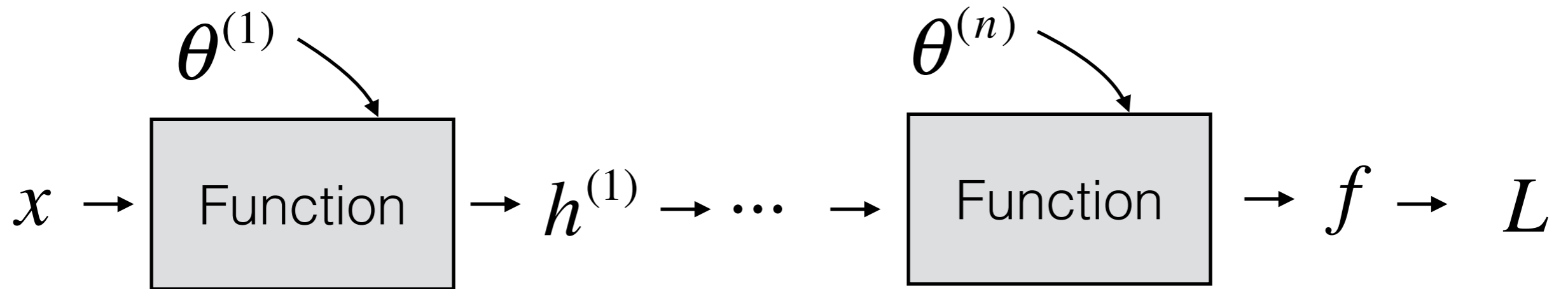


Backward Propagation:

L

Recap

Forward Propagation:

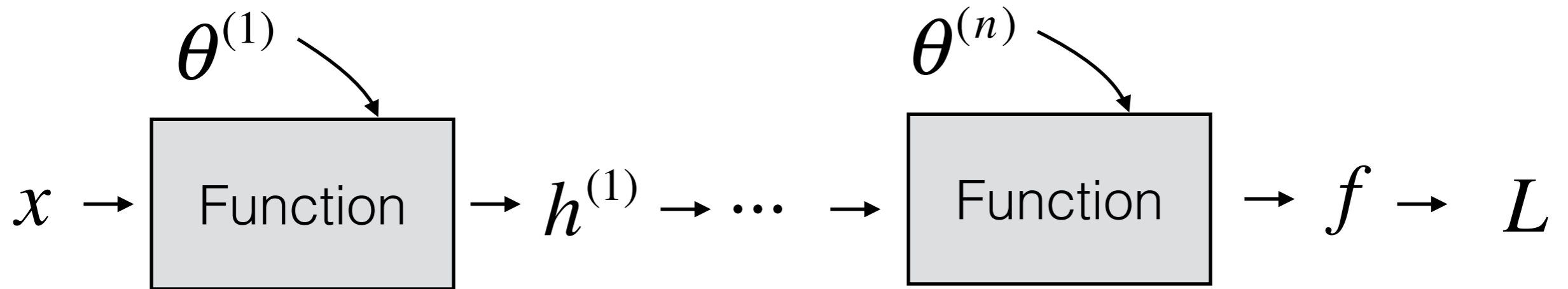


Backward Propagation:

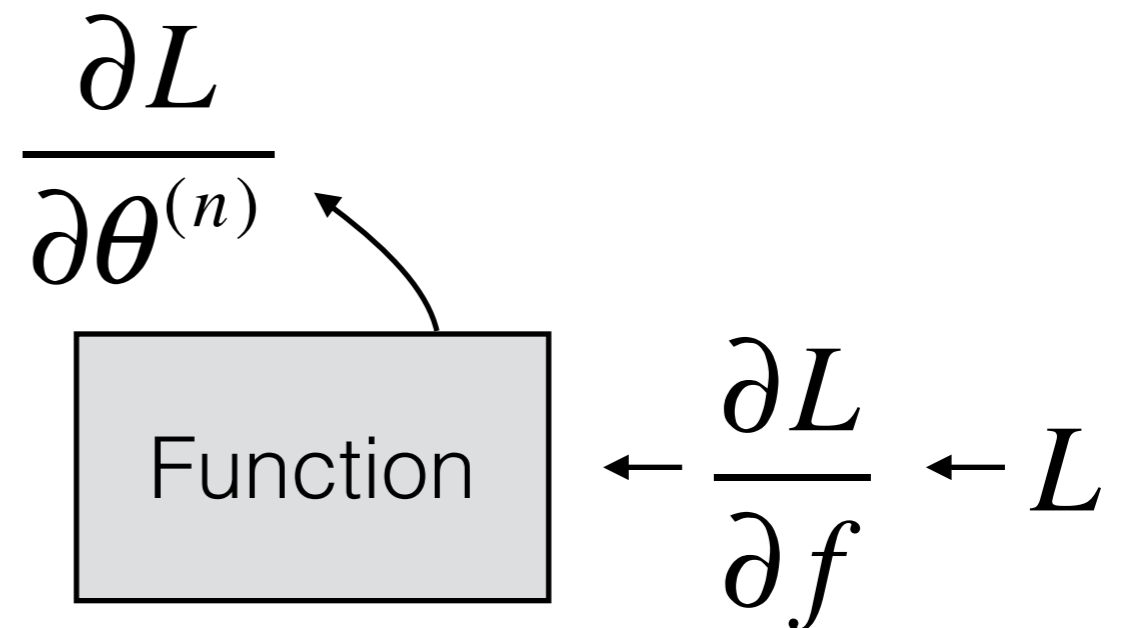
$$\frac{\partial L}{\partial f} \leftarrow L$$

Recap

Forward Propagation:

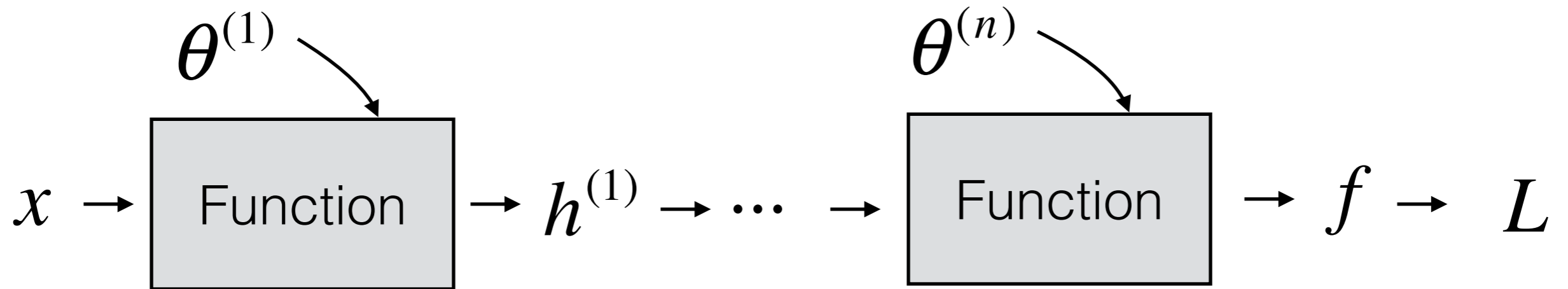


Backward Propagation:

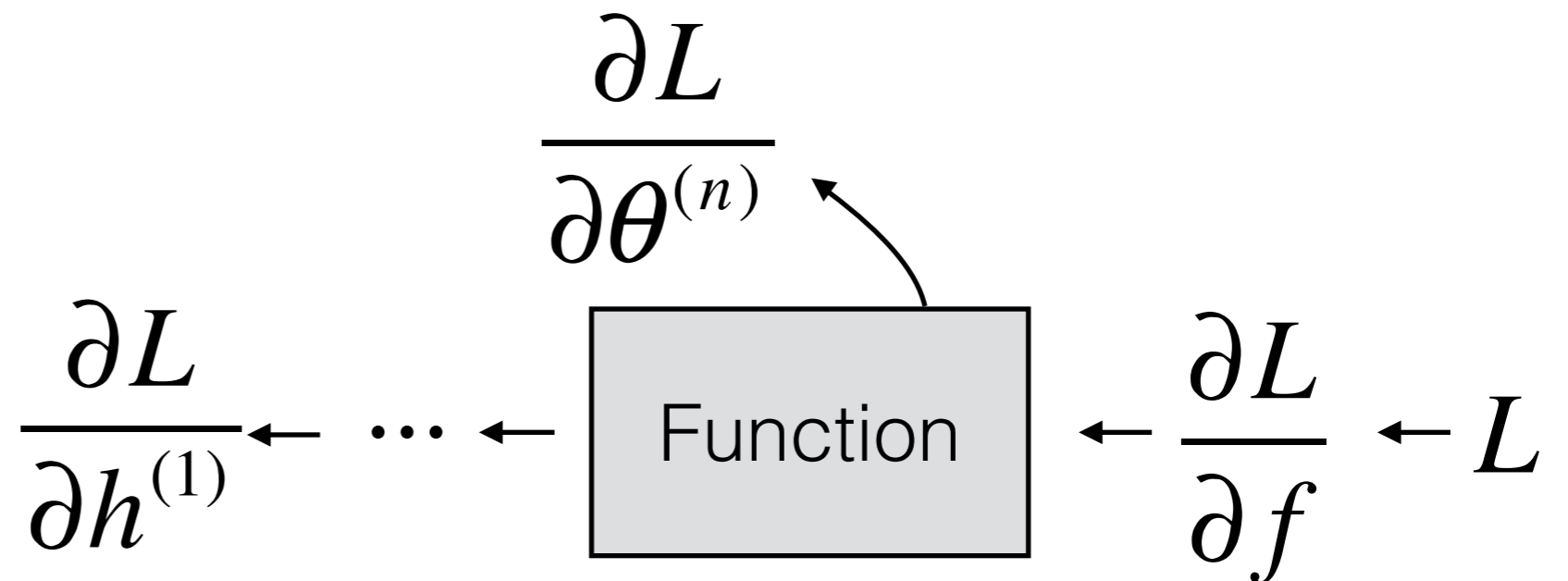


Recap

Forward Propagation:

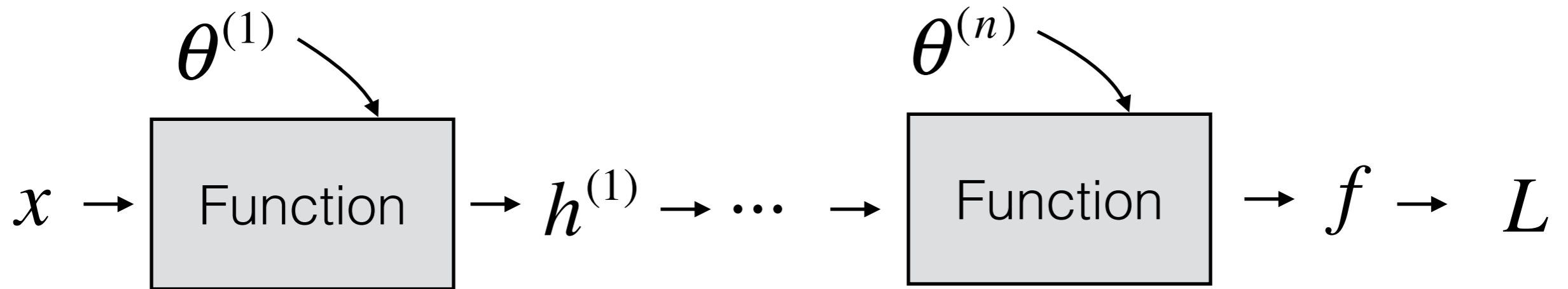


Backward Propagation:

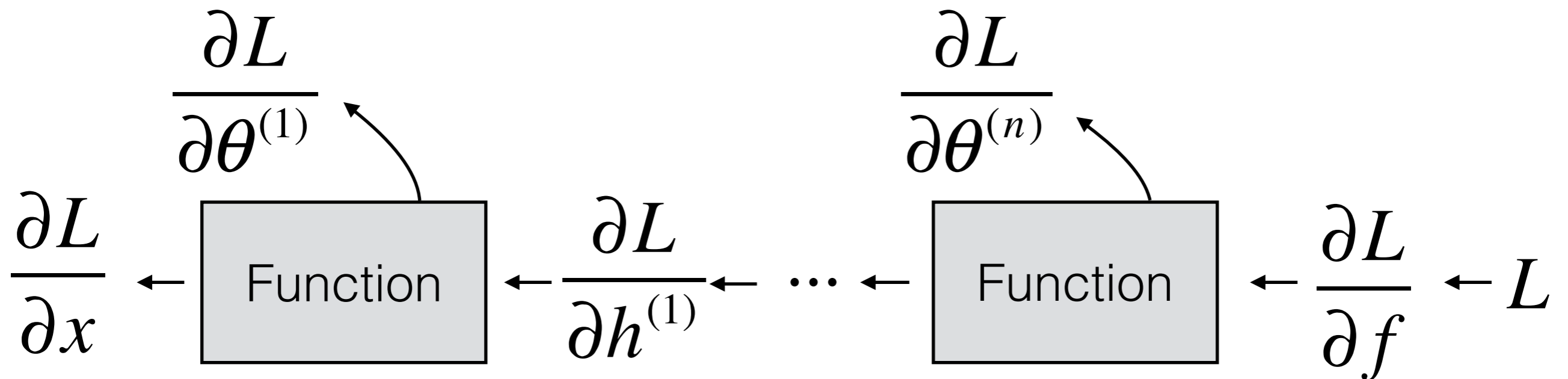


Recap

Forward Propagation:



Backward Propagation:



Questions?

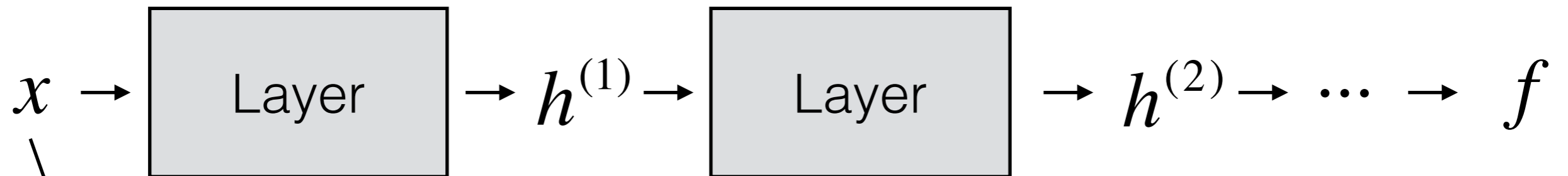
30s cat picture break



CNNs

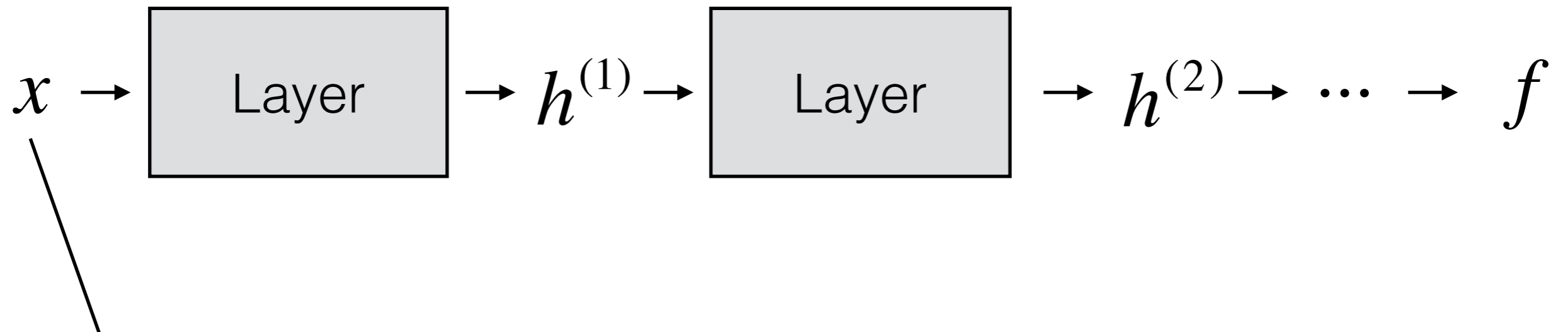
It's just neural networks
with 3D activations

What shape should the activations have?



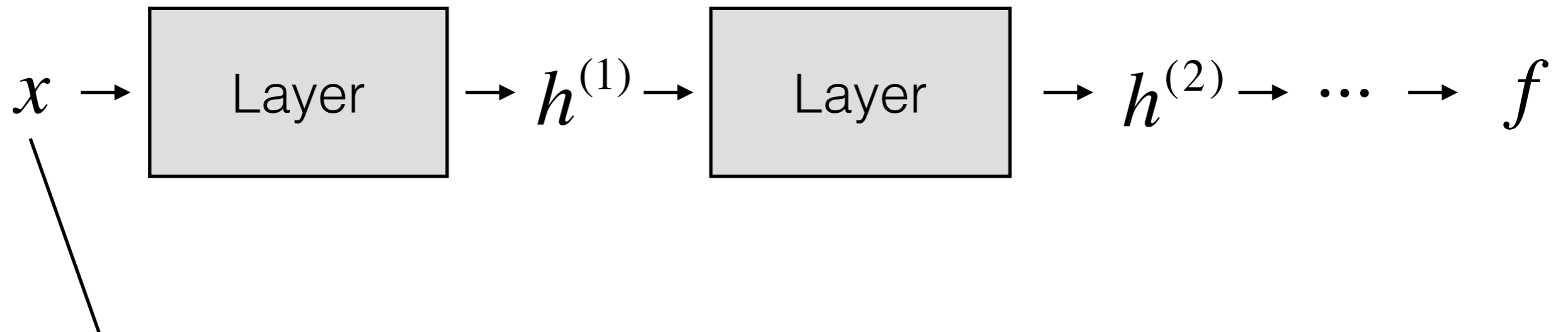
- The input is an image, which is 3D (RGB channel, height, width)

What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure

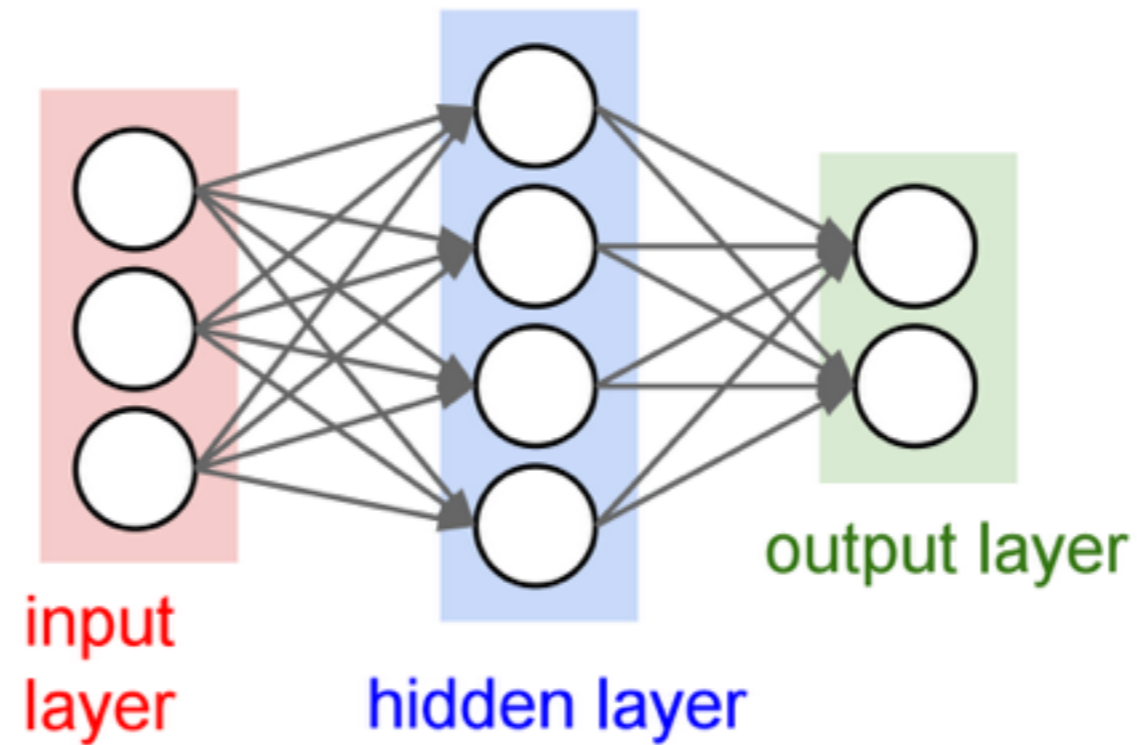
What shape should the activations have?



- The input is an image, which is 3D (RGB channel, height, width)
- We could flatten it to a 1D vector, but then we lose structure
- What about keeping everything in 3D?

3D Activations

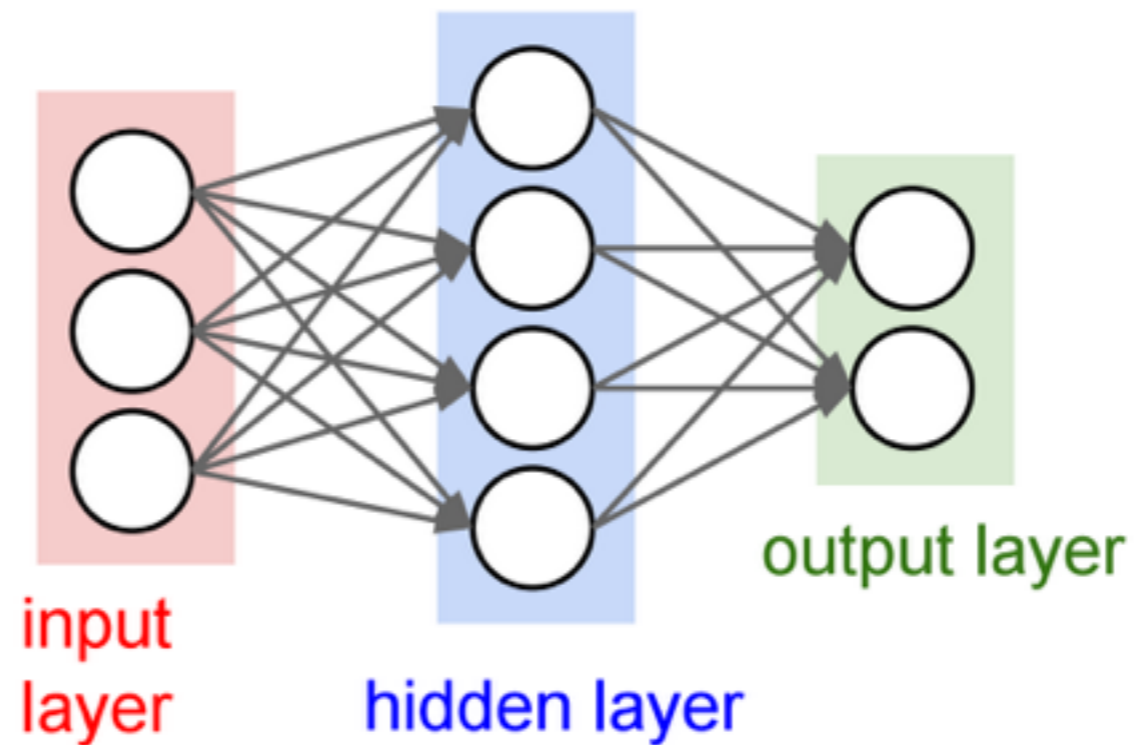
before:



(1D vectors)

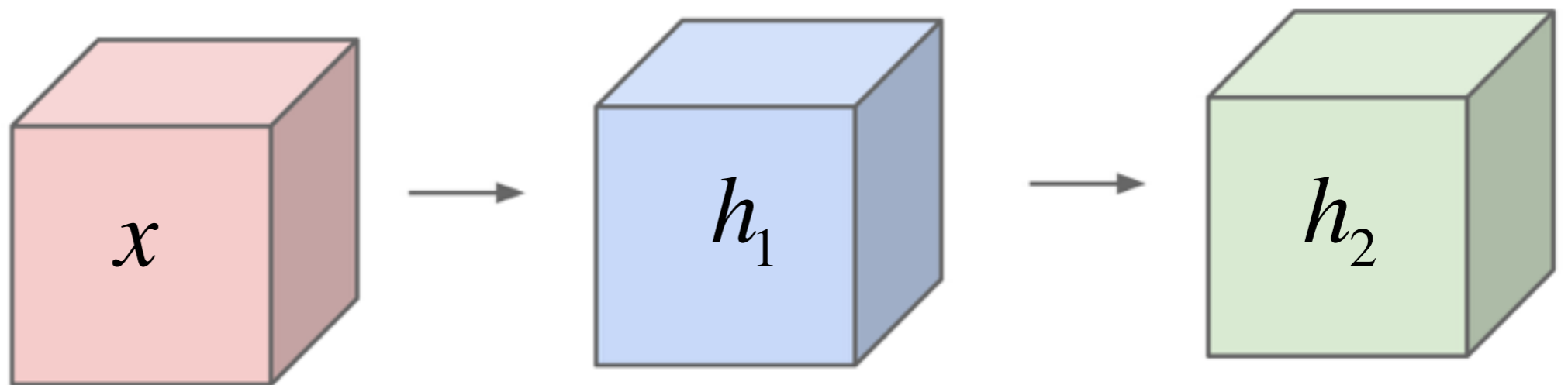
3D Activations

before:



(1D vectors)

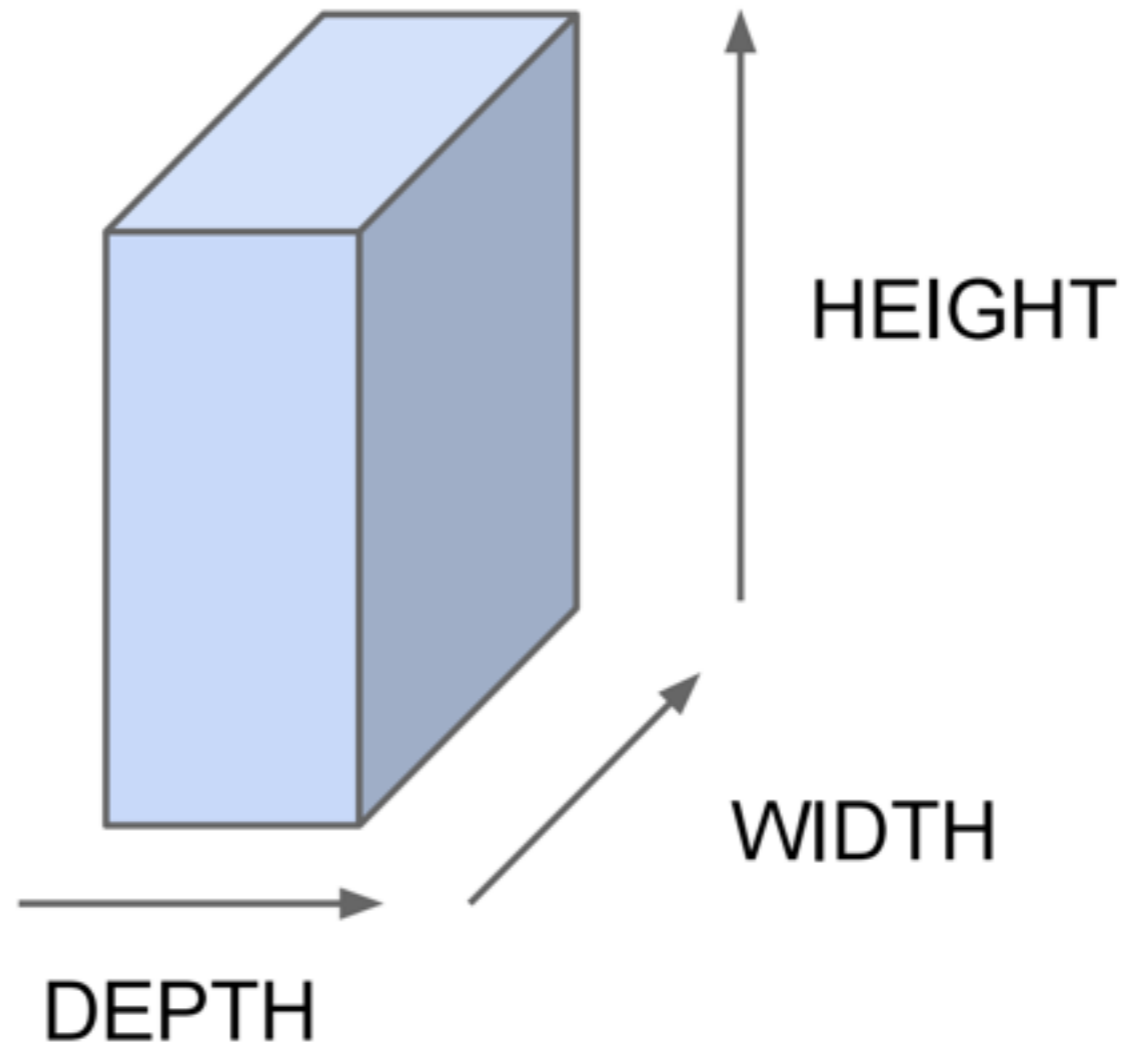
now:



(3D arrays)

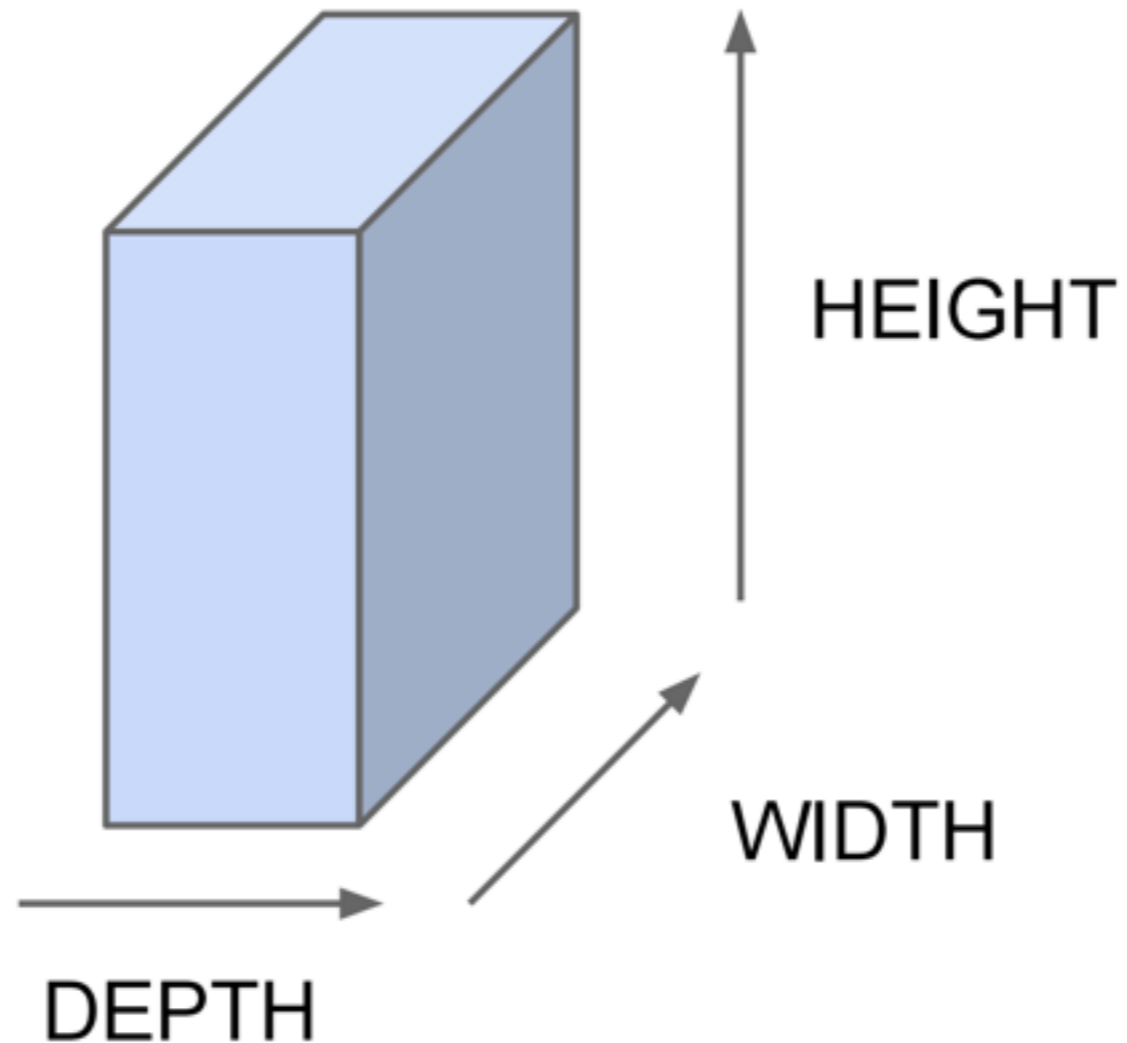
3D Activations

All Neural Net activations arranged in **3 dimensions**:



3D Activations

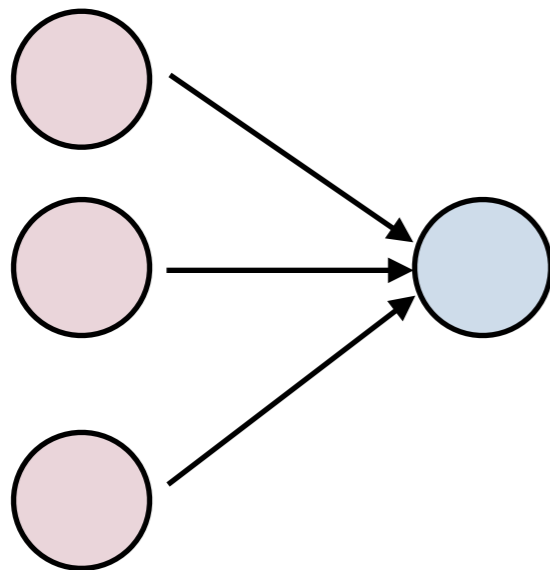
All Neural Net activations arranged in **3 dimensions**:



For example, a CIFAR-10 image is a $3 \times 32 \times 32$ volume (3 depth — RGB channels, 32 height, 32 width)

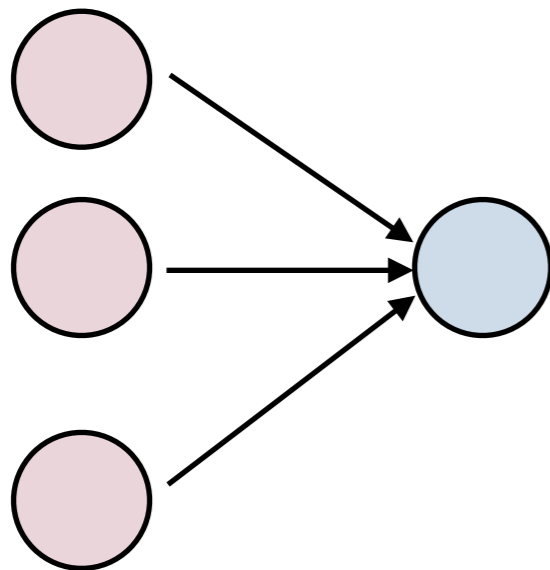
3D Activations

1D Activations:

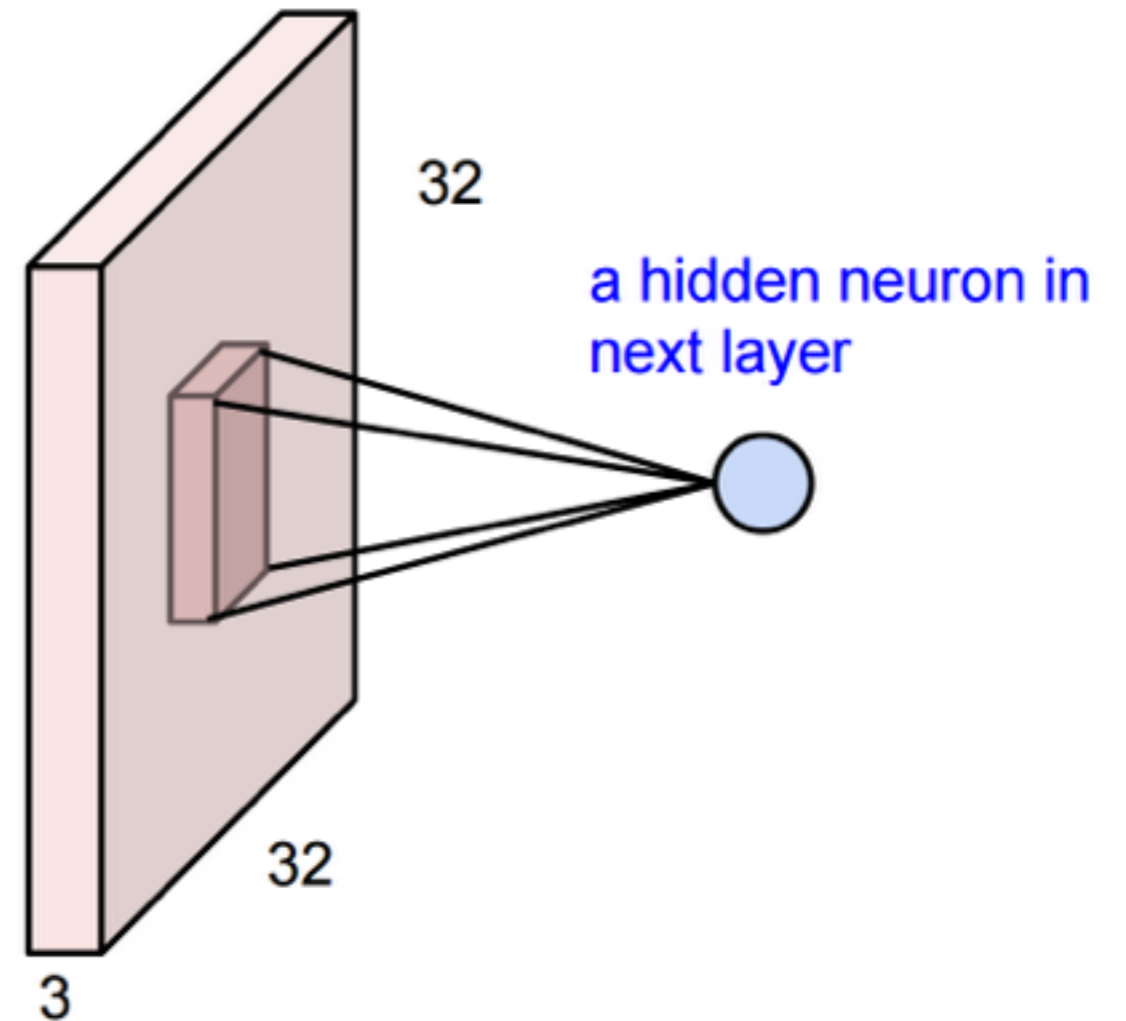


3D Activations

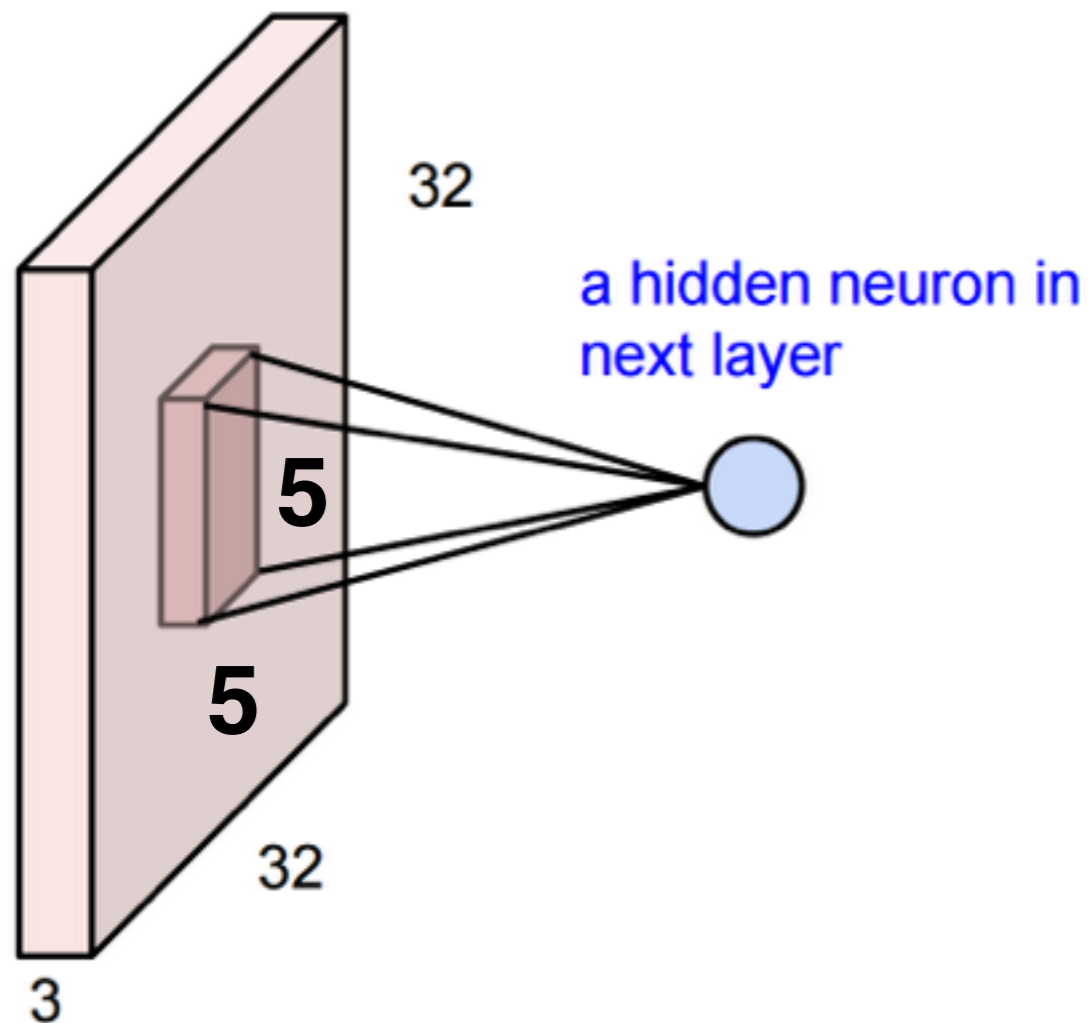
1D Activations:



3D Activations:

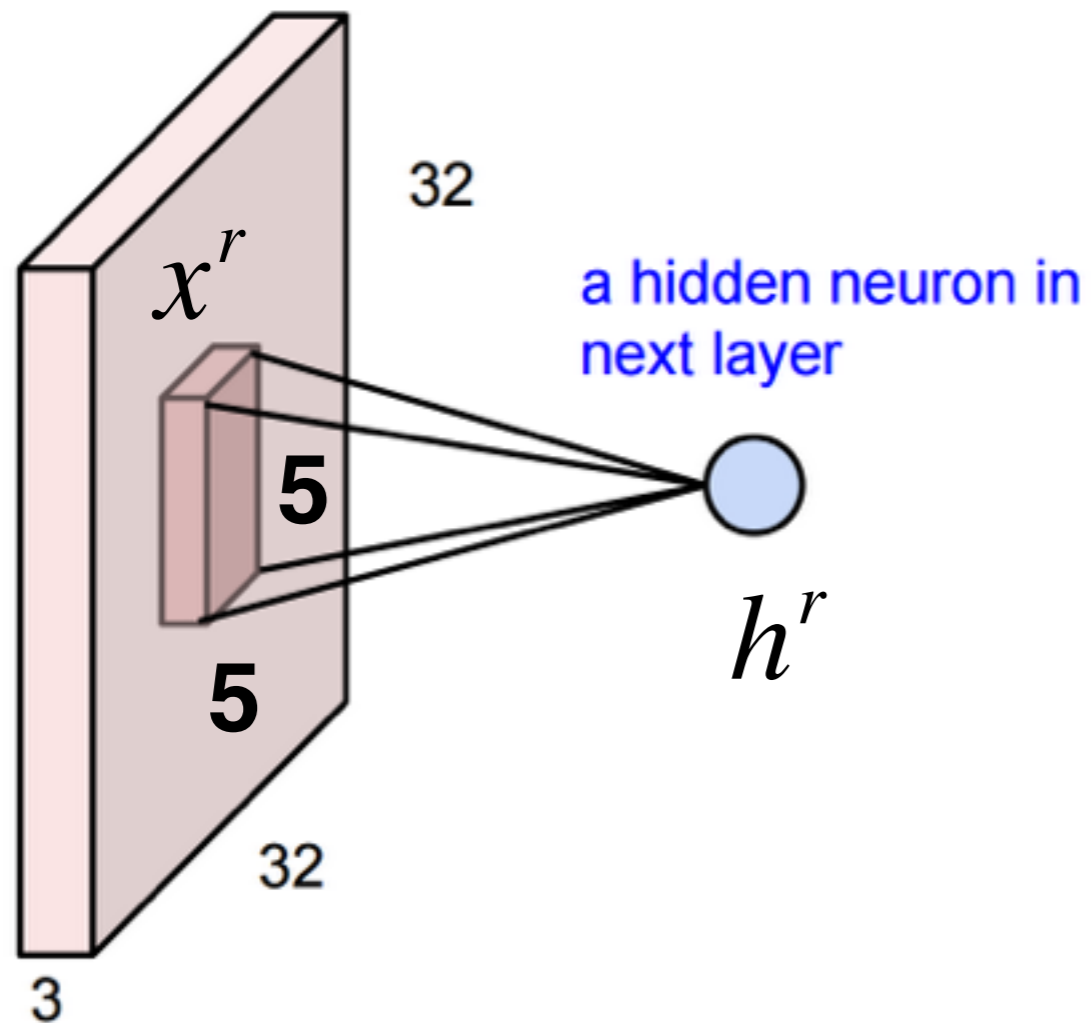


3D Activations



- The input is $3 \times 32 \times 32$
- This neuron depends on a $3 \times 5 \times 5$ chunk of the input
- The neuron also has a $3 \times 5 \times 5$ set of weights and a bias (scalar)

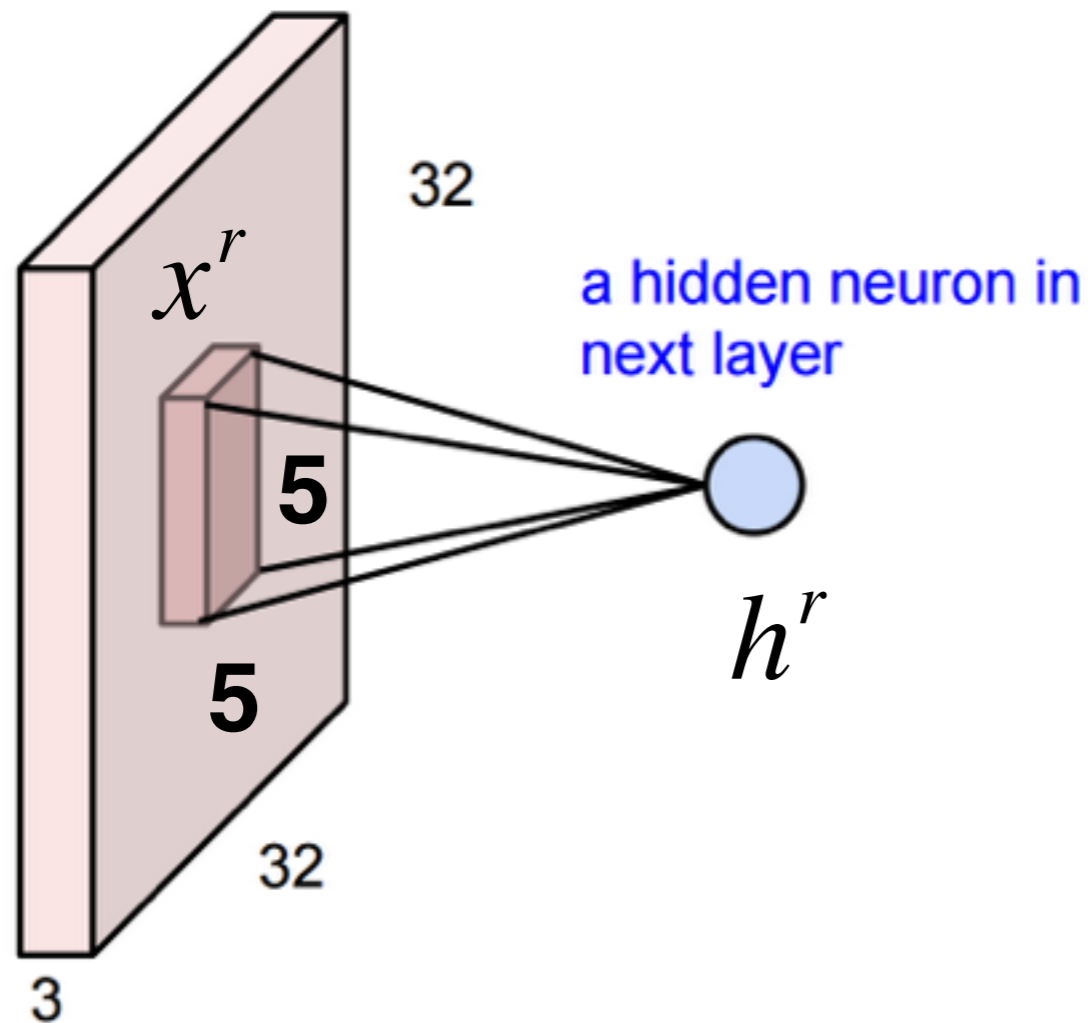
3D Activations



Example: consider the region of the input " x^r "

With output neuron h^r

3D Activations



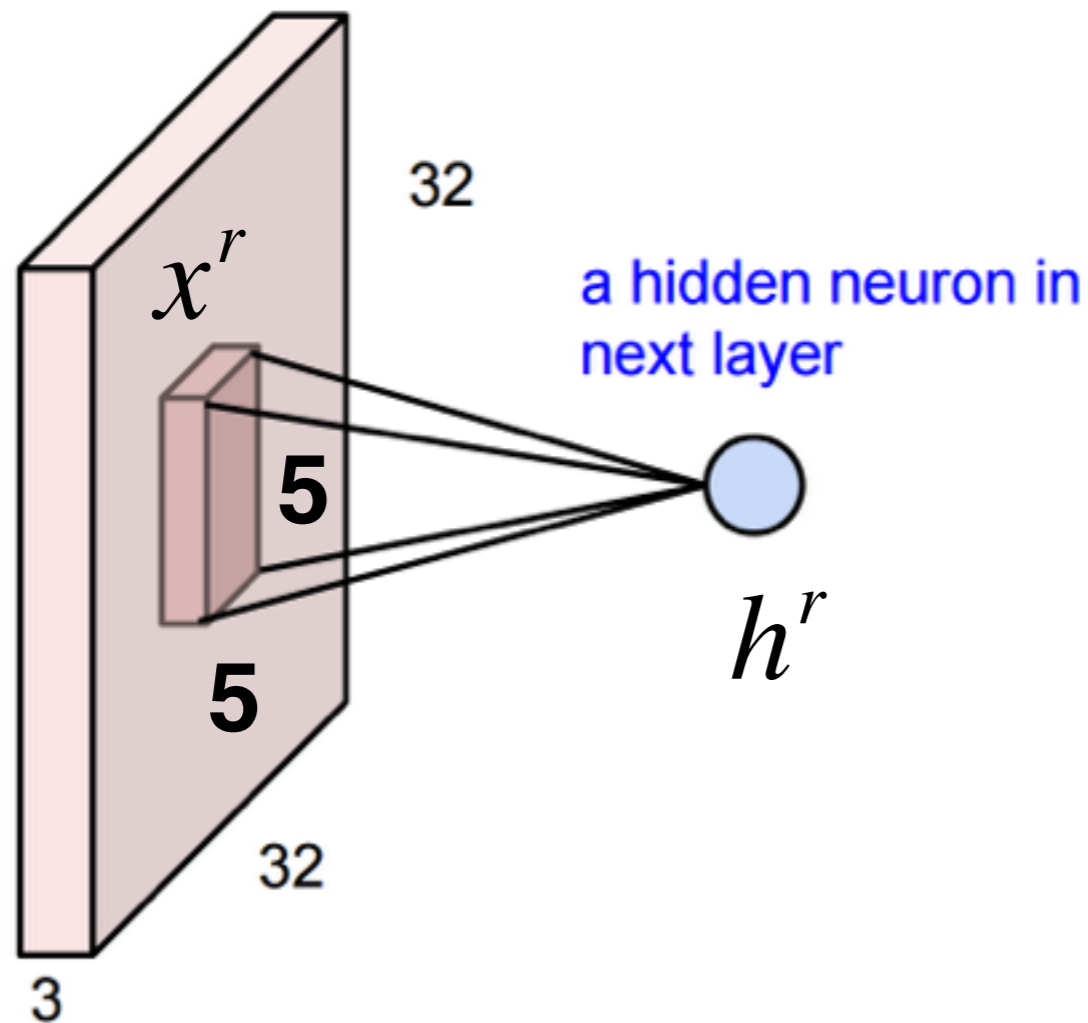
Example: consider the region of the input “ x^r ”

With output neuron h^r

Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

3D Activations



Example: consider the region of the input “ x^r ”

With output neuron h^r

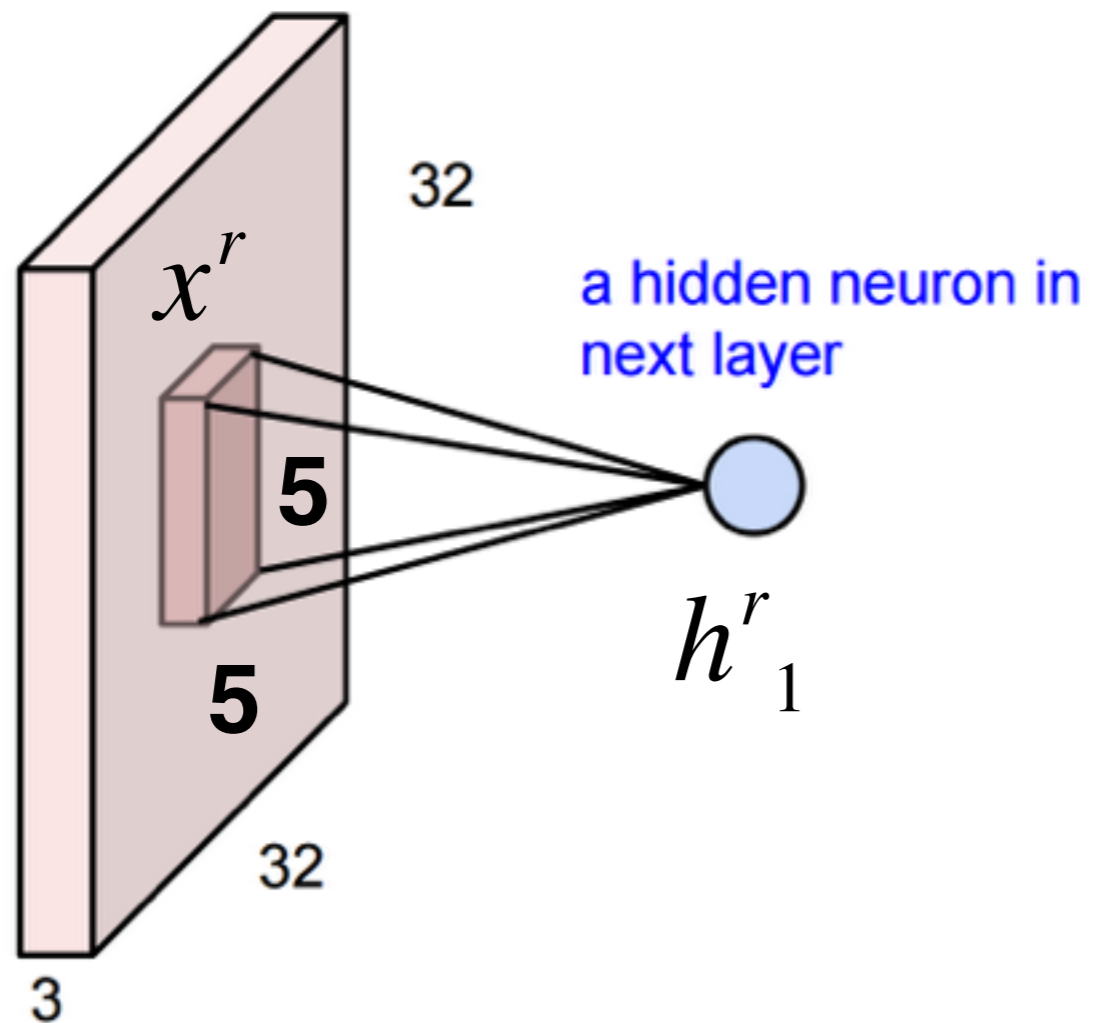
Then the output is:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

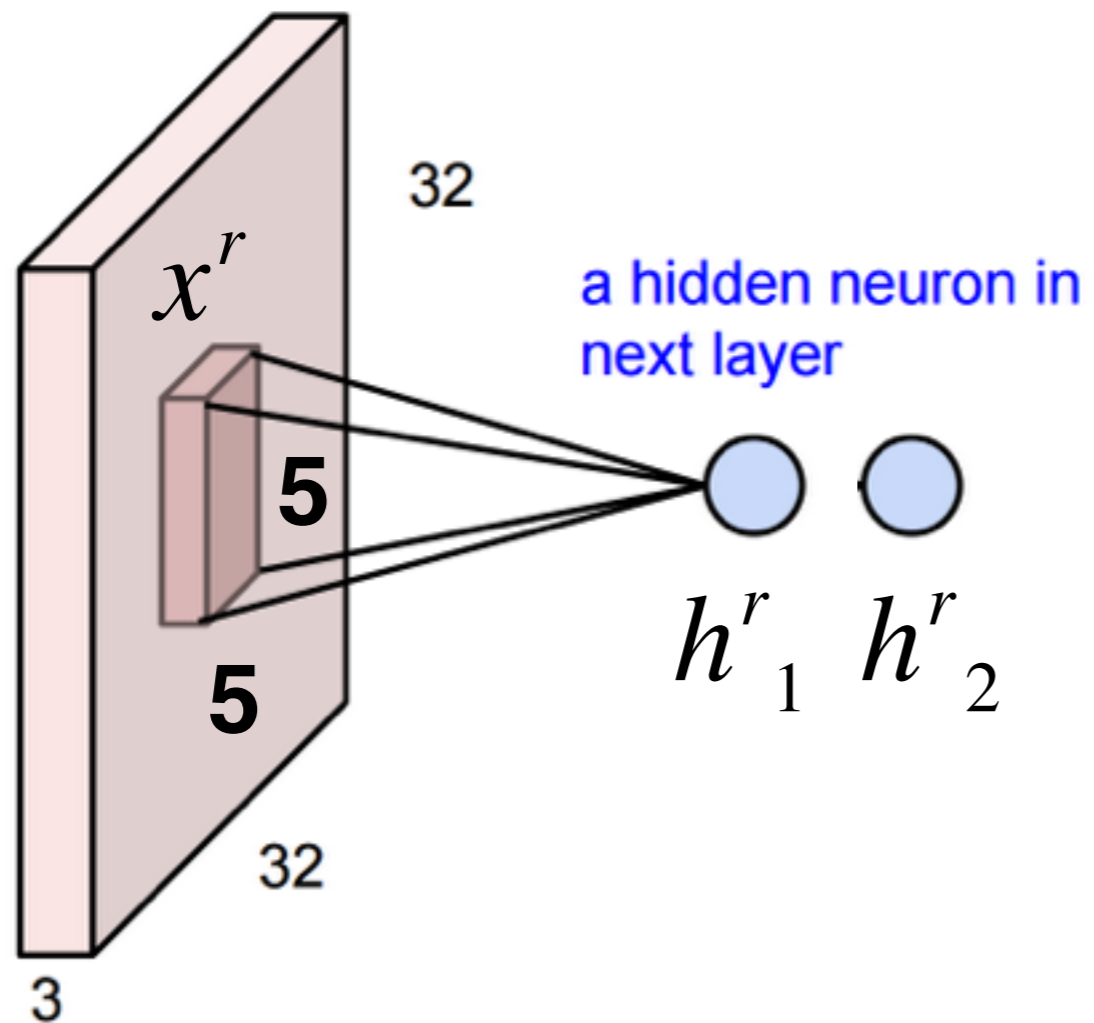


Sum over 3 axes

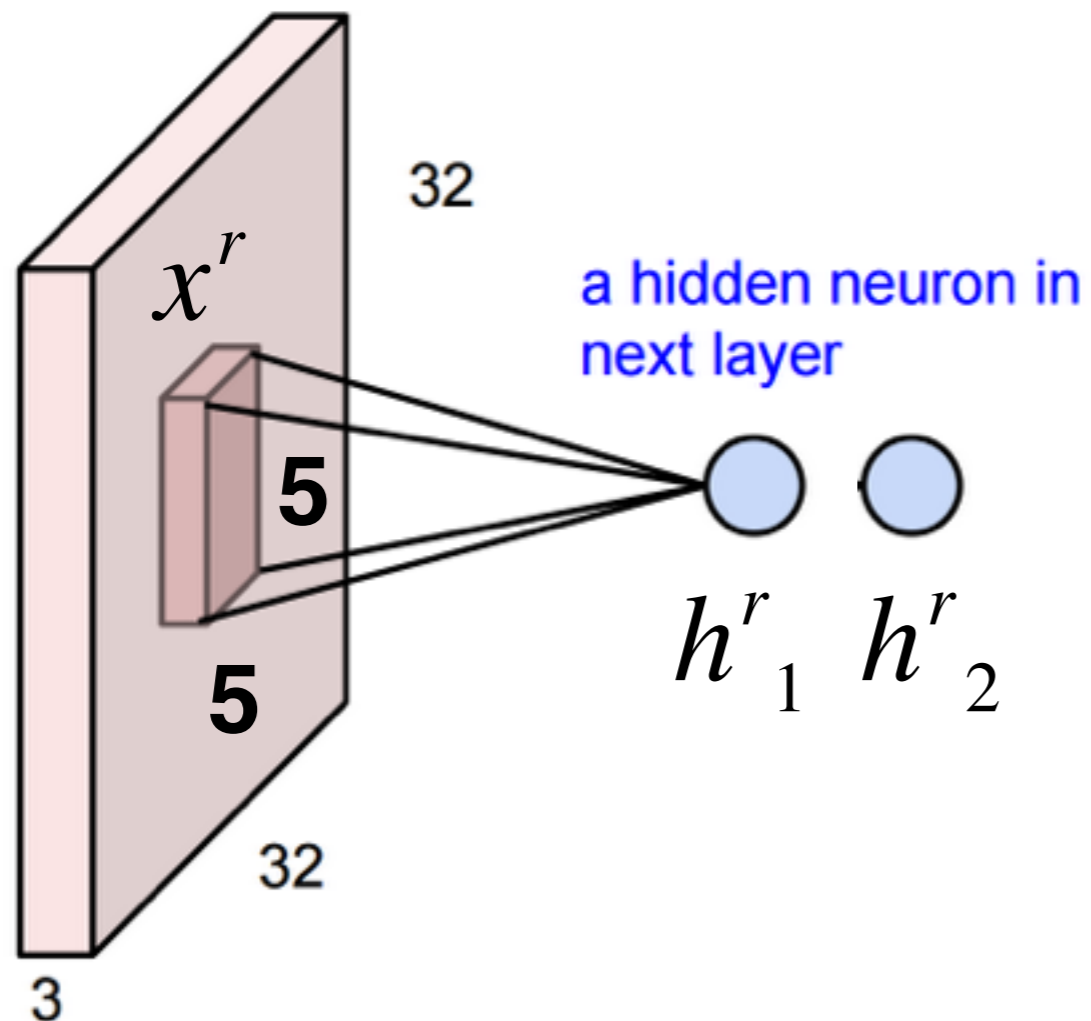
3D Activations



3D Activations



3D Activations

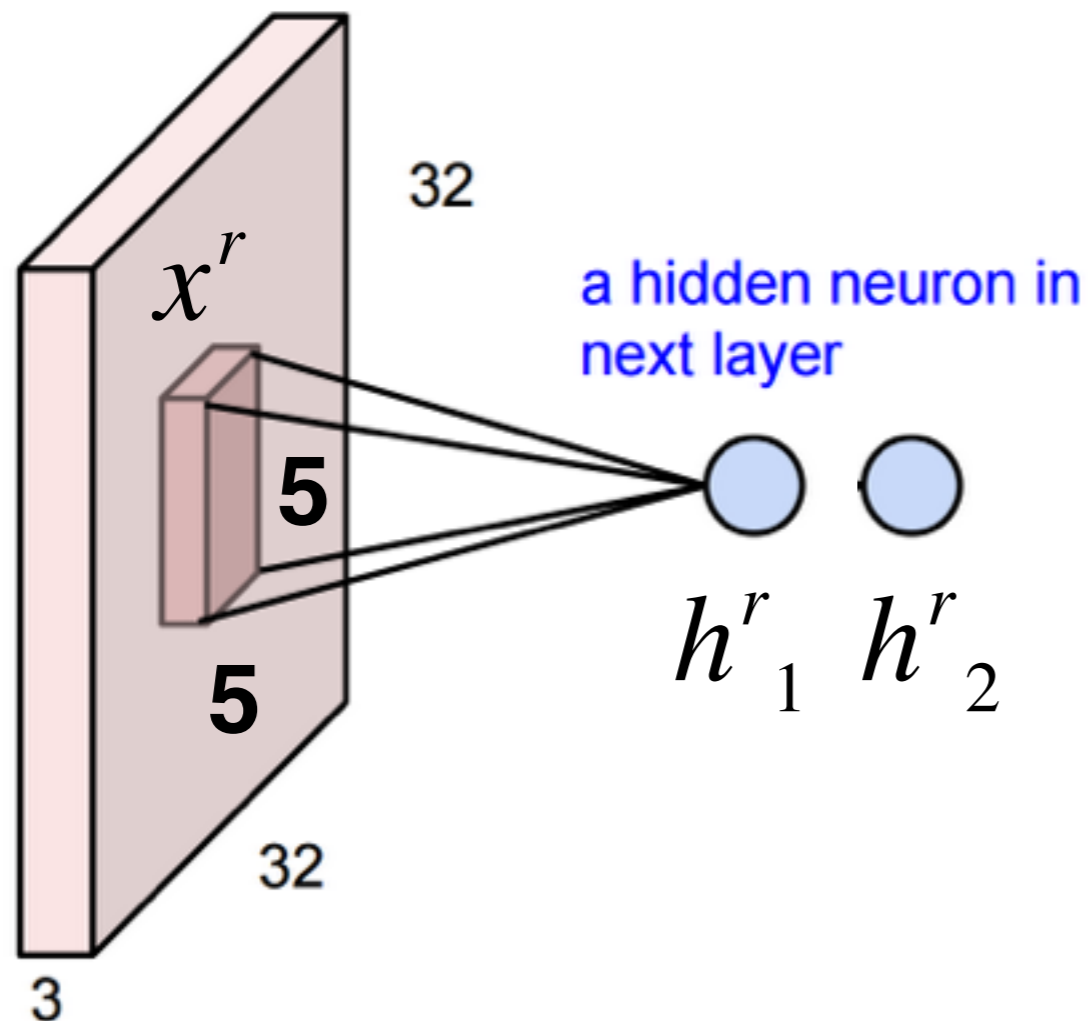


With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_1$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_2$$

3D Activations



With **2** output neurons

$$h^r_1 = \sum_{ijk} x^r_{ijk} W_{1ijk} + b_{1}$$

$$h^r_2 = \sum_{ijk} x^r_{ijk} W_{2ijk} + b_{2}$$

3D Activations

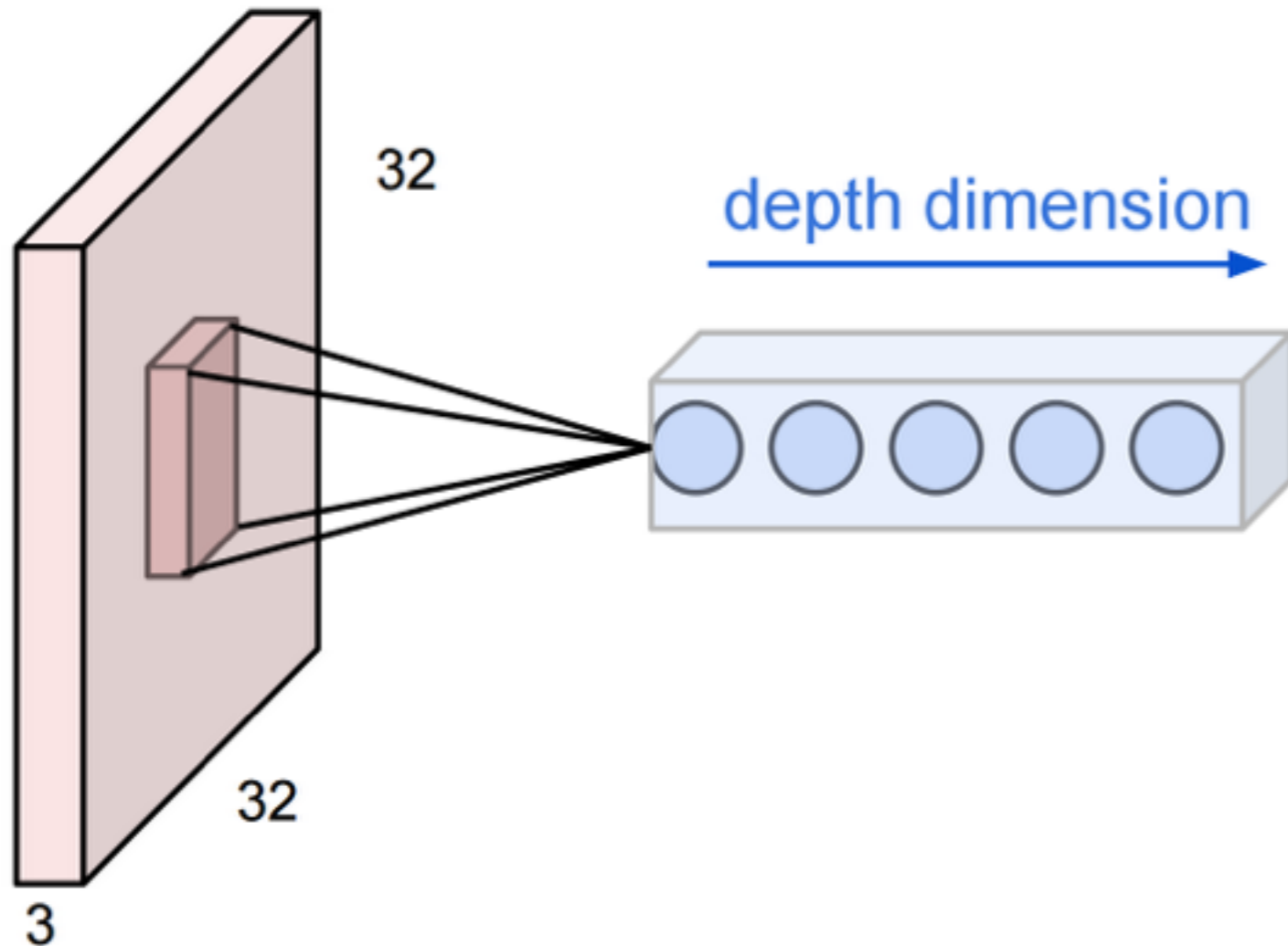
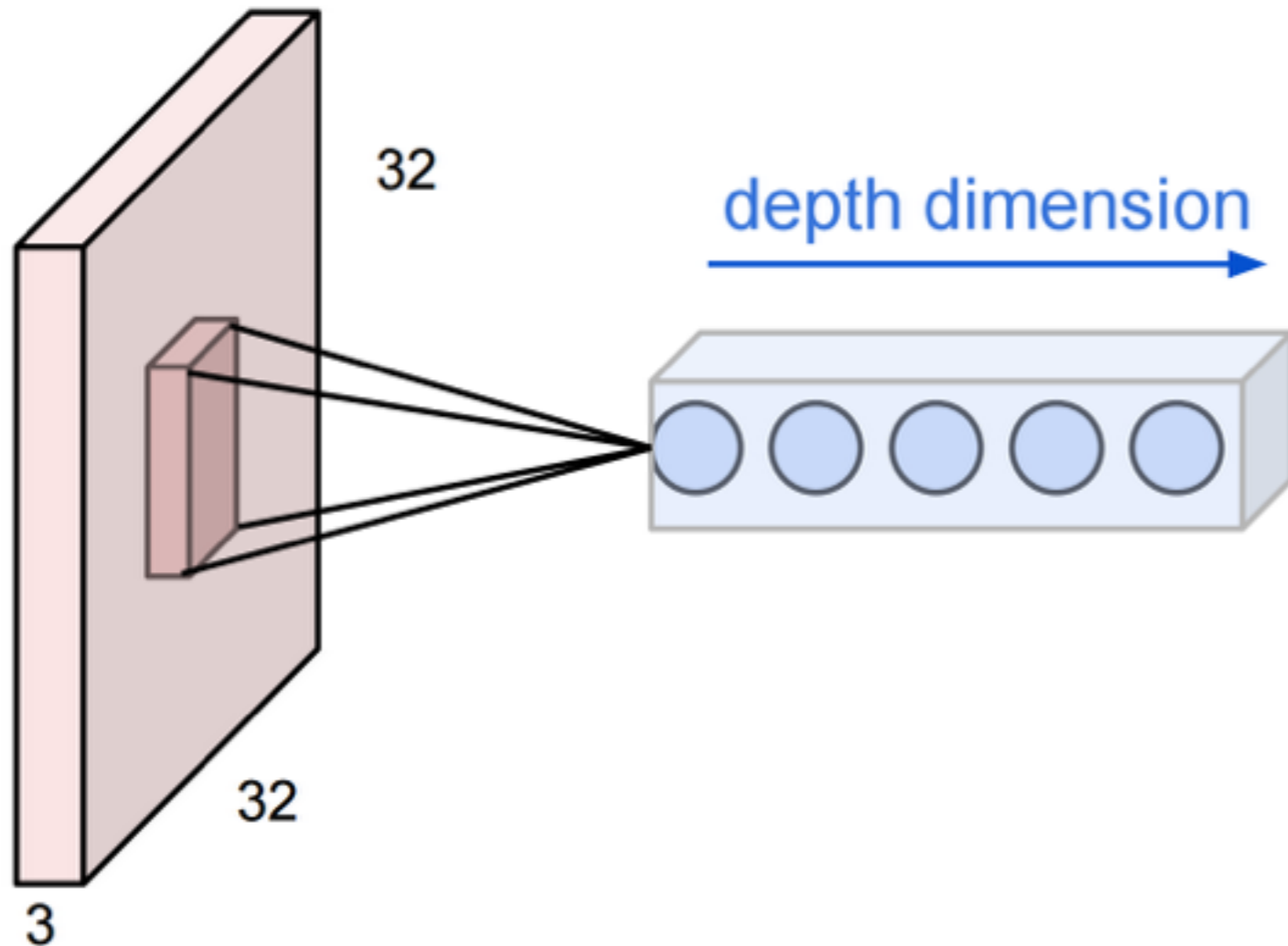


Figure: Andrej Karpathy

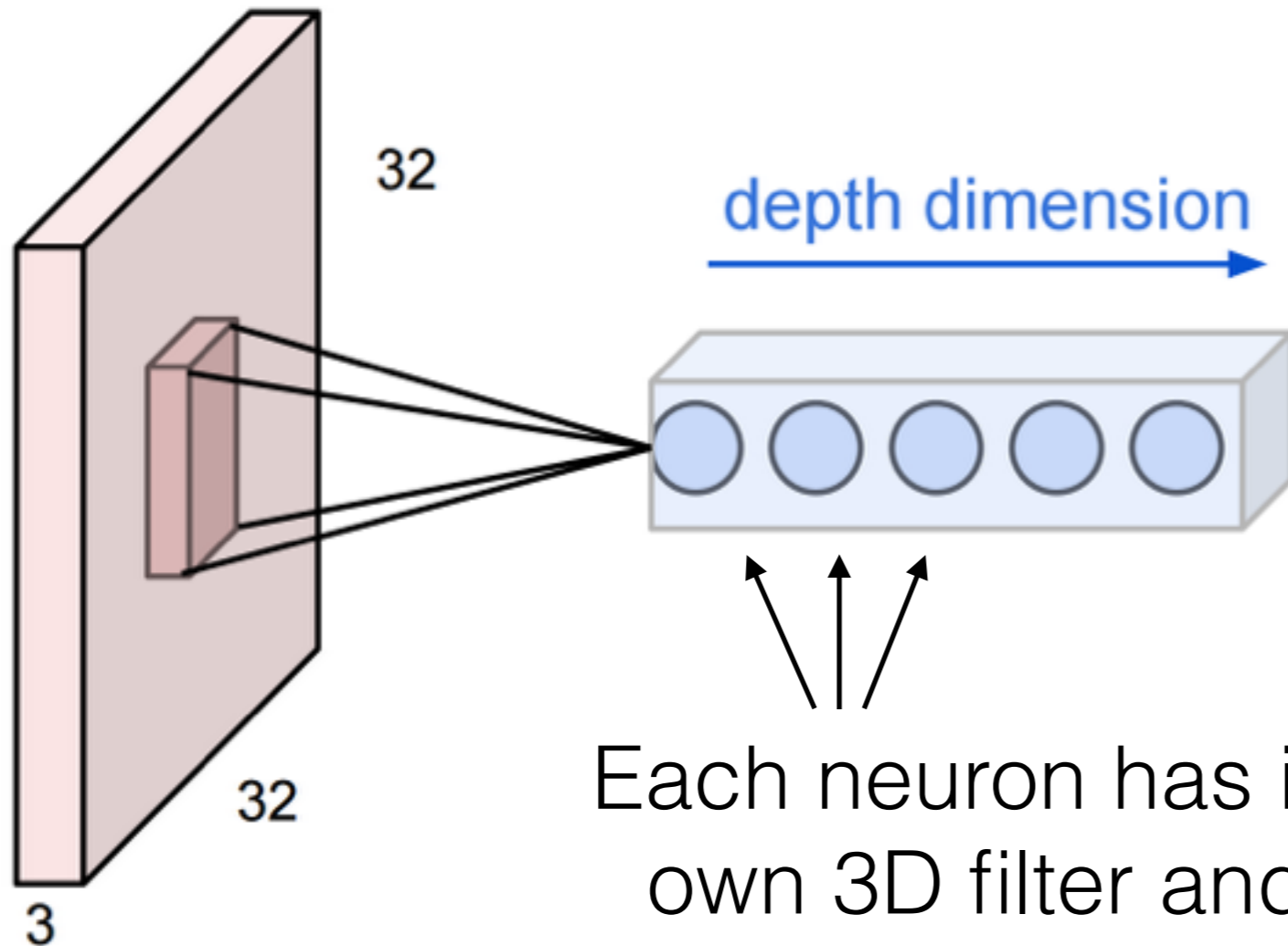
3D Activations



We can keep adding more outputs

These form a column in the output volume:
[depth x 1 x 1]

3D Activations

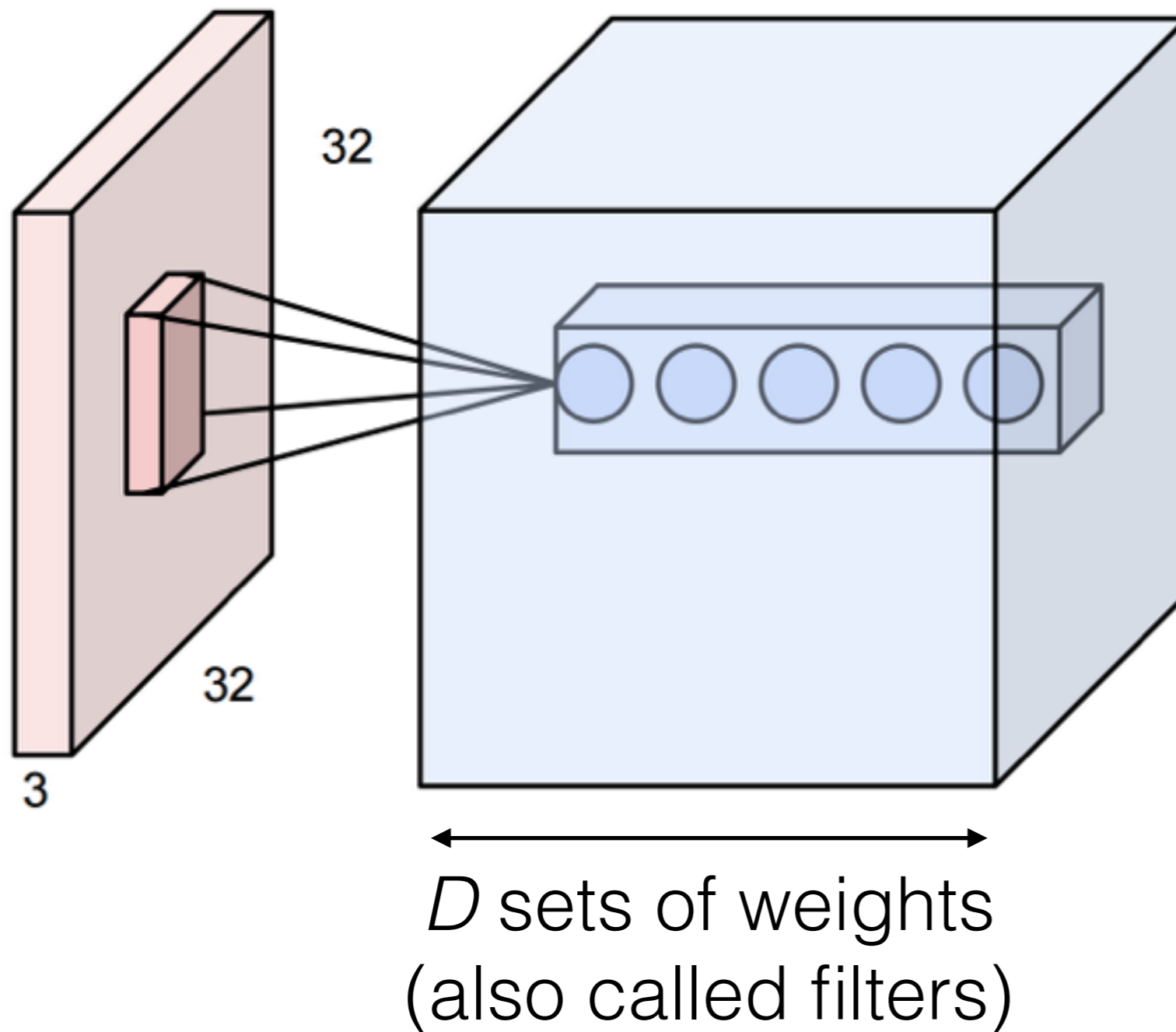


Each neuron has its own 3D filter and own (scalar) bias

We can keep adding more outputs

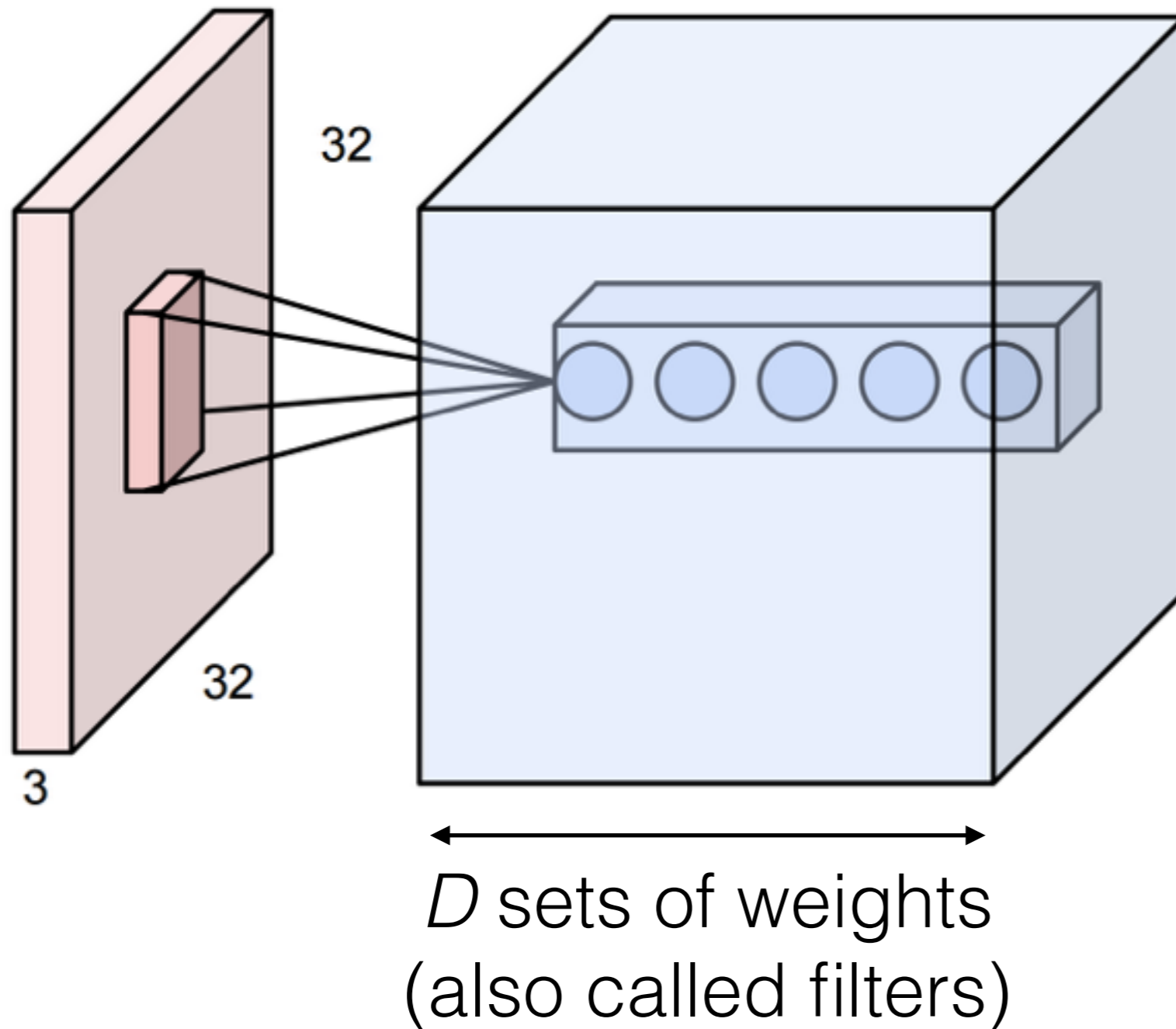
These form a column in the output volume:
[depth x 1 x 1]

3D Activations



Now repeat this
across the input

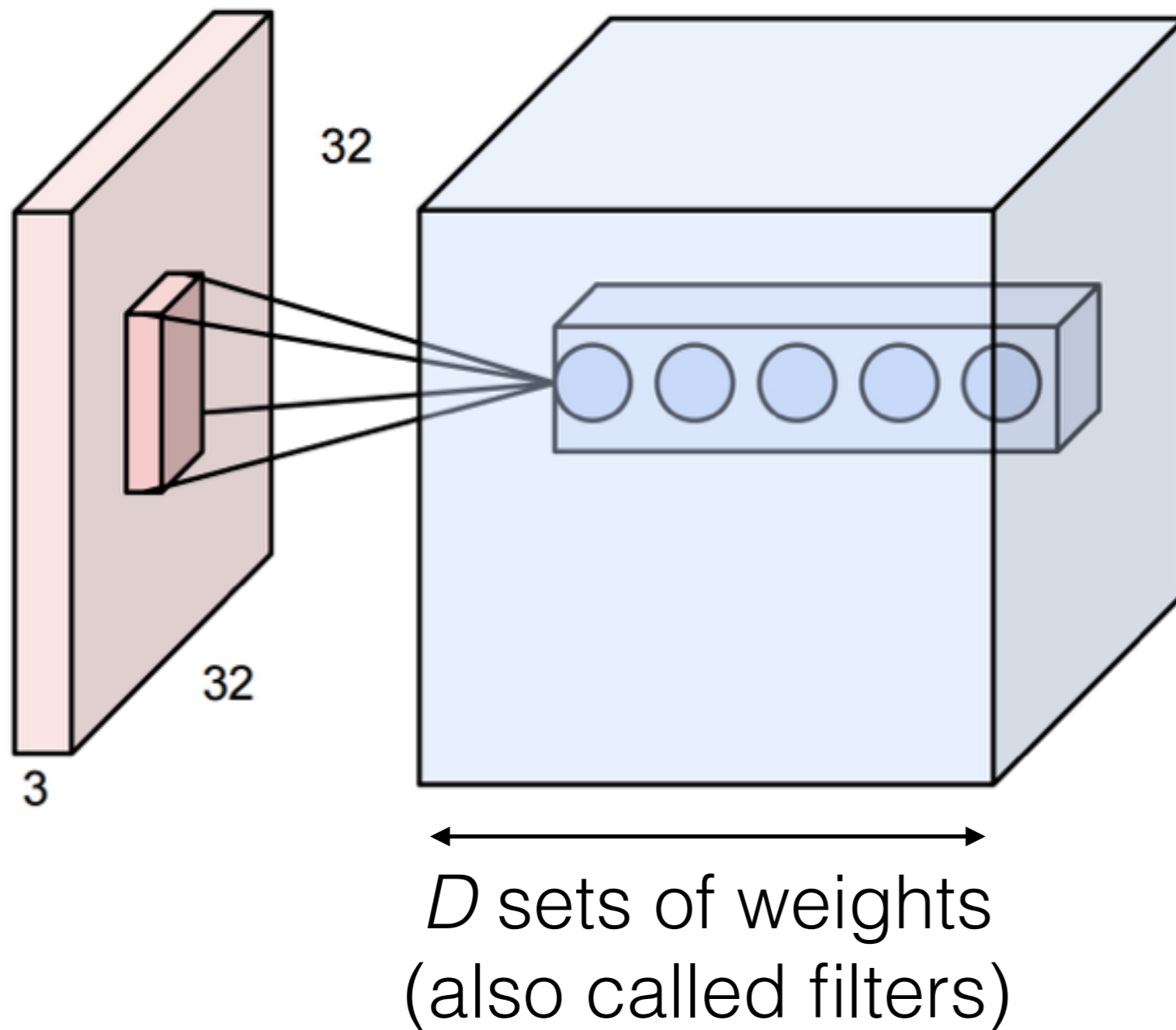
3D Activations



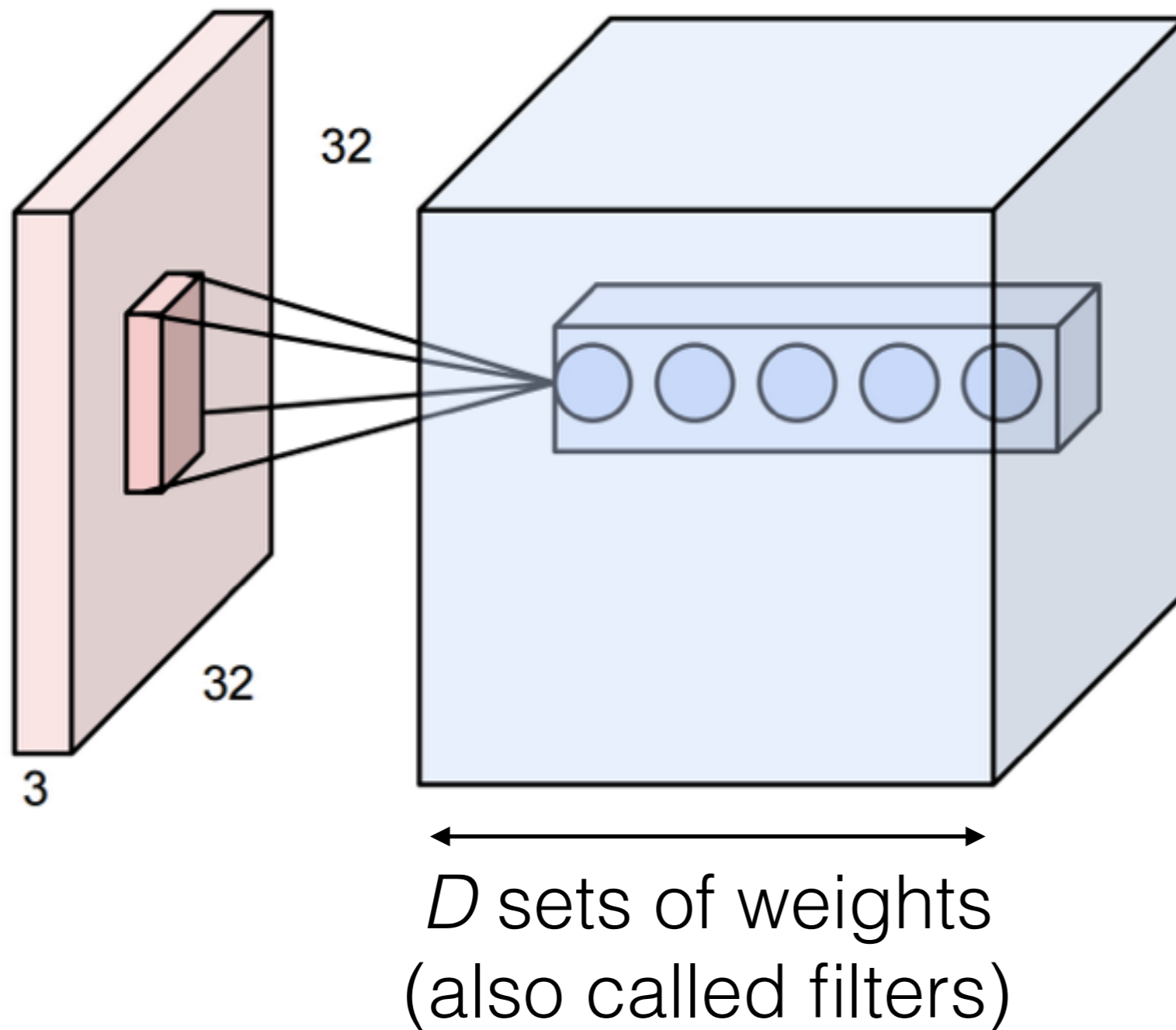
Now repeat this
across the input

Weight sharing:
Each filter shares
the same weights
(but each depth
index has its own
set of weights)

3D Activations

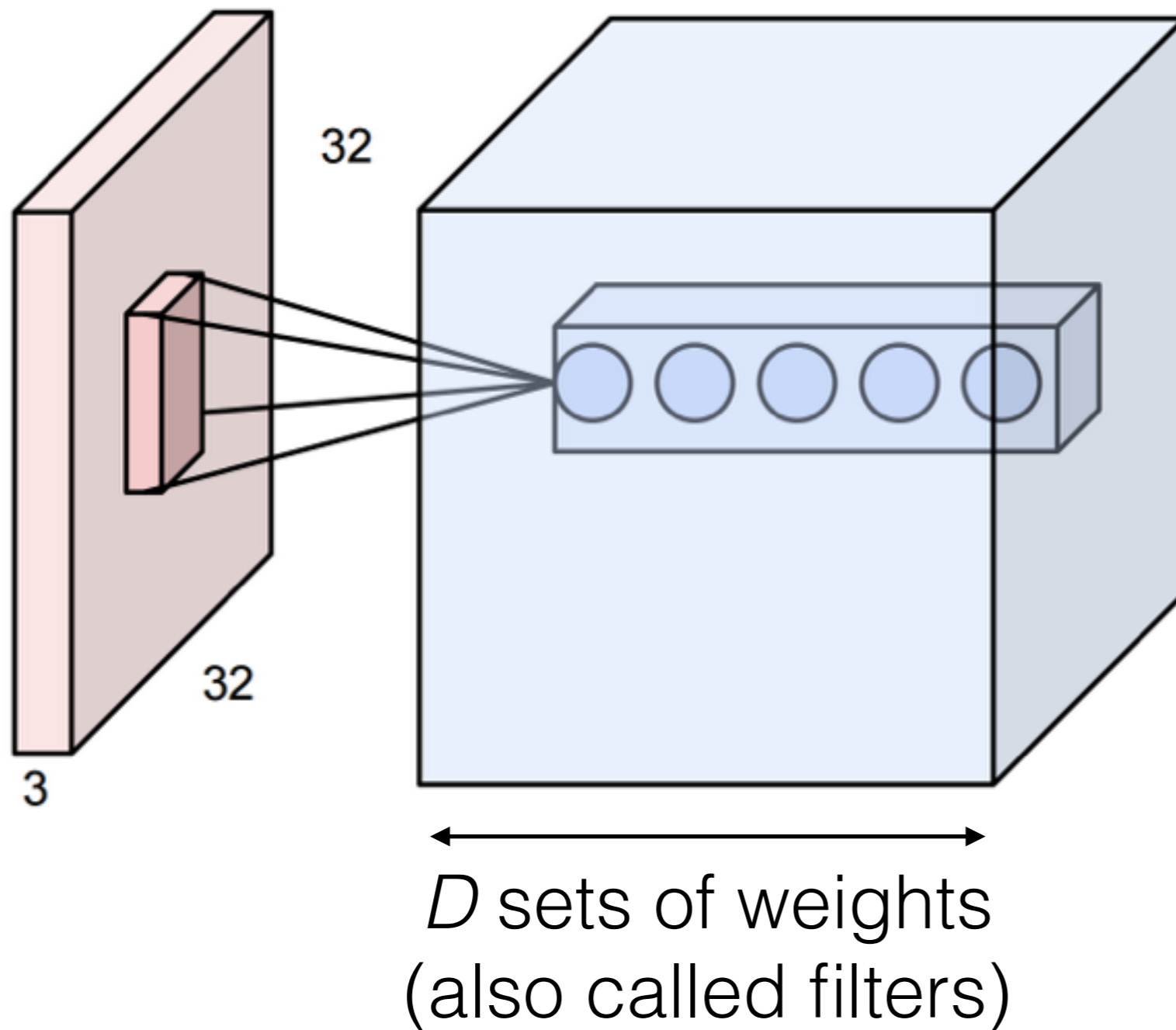


3D Activations



With weight sharing,
this is called
convolution

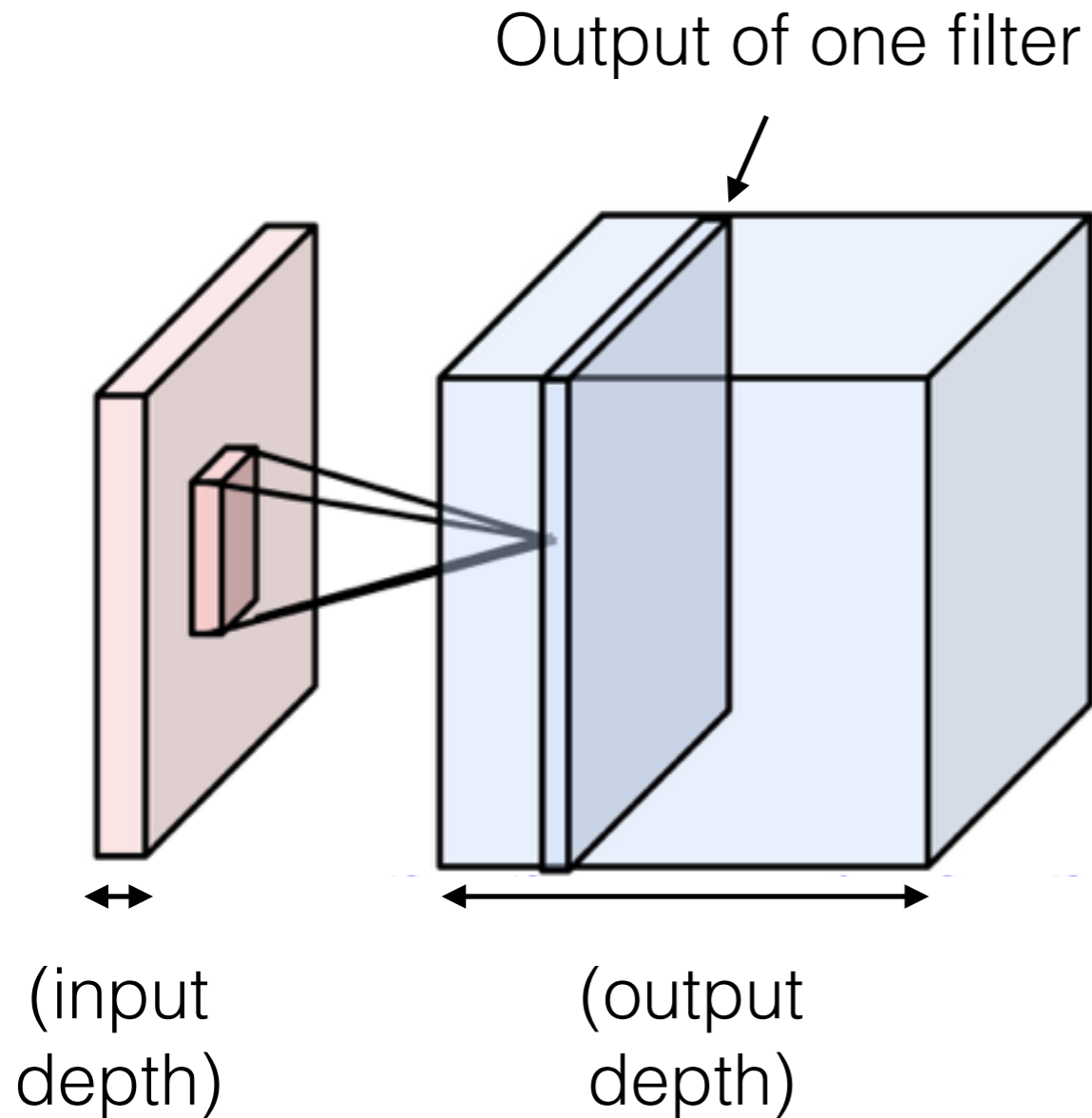
3D Activations



With weight sharing,
this is called
convolution

Without weight sharing,
this is called a
locally connected layer

3D Activations

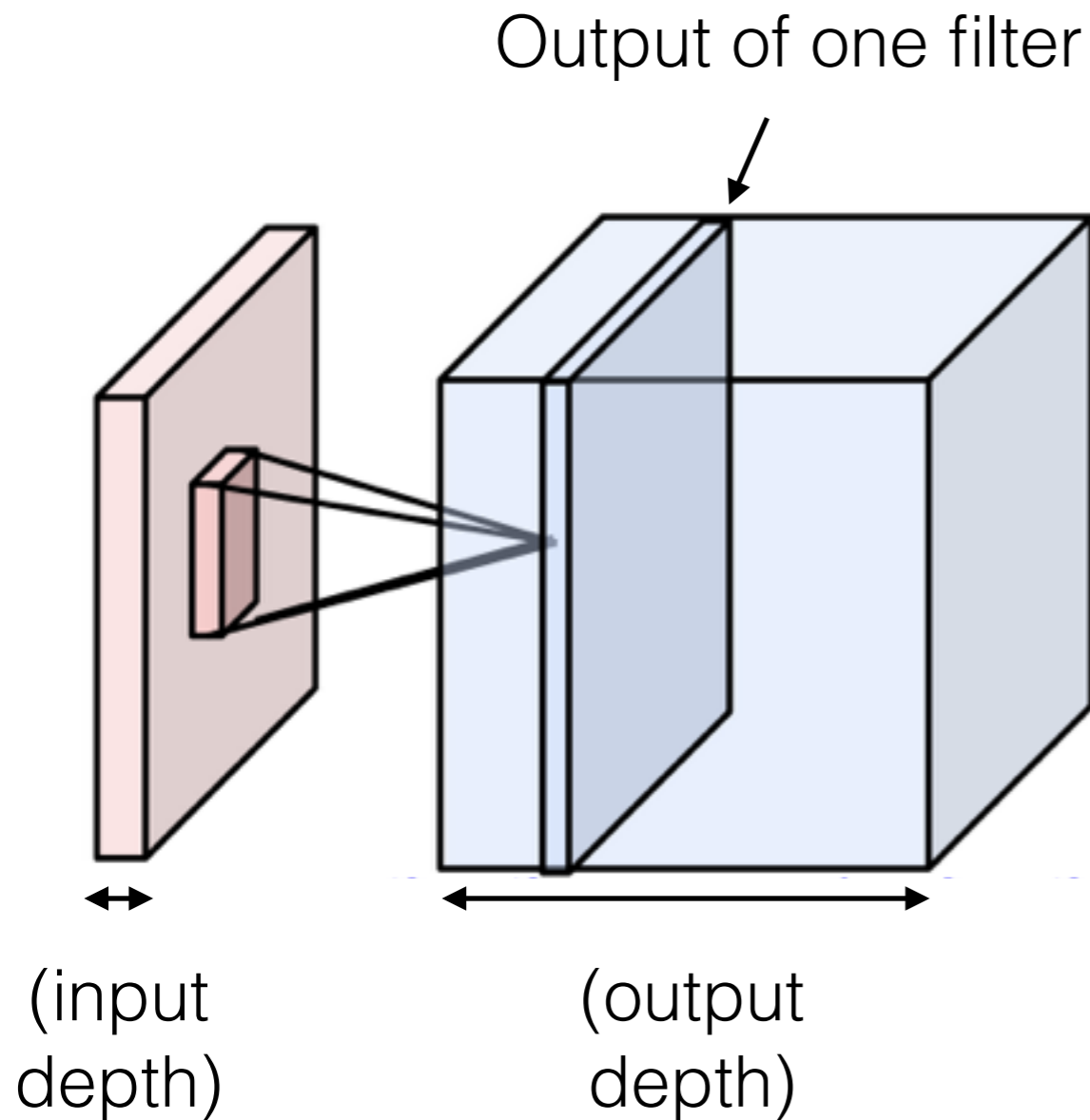


One set of weights gives one slice in the output

To get a 3D output of depth D , use D different filters

In practice, CNNs use many filters (~ 64 to 1024)

3D Activations



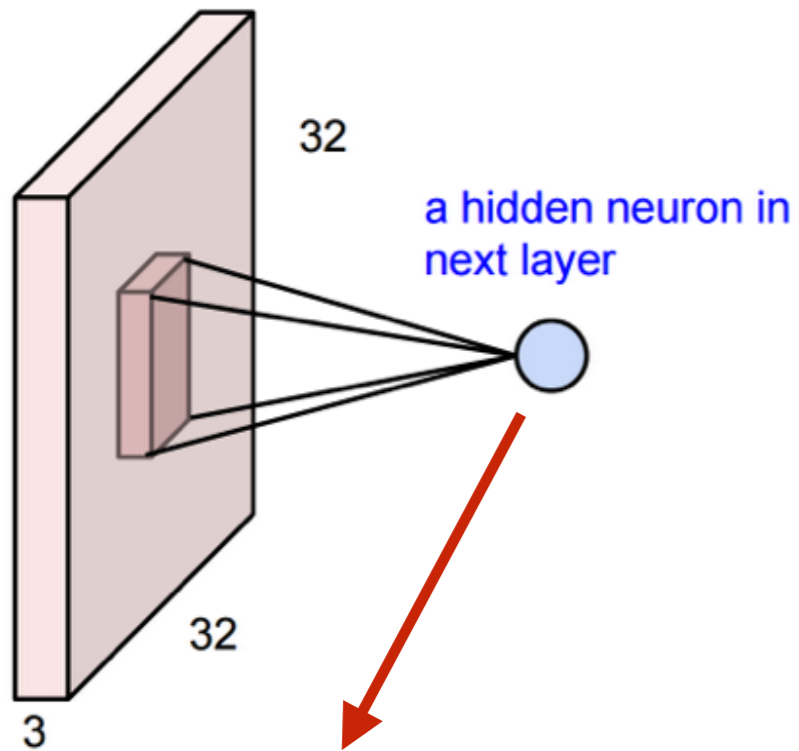
One set of weights gives one slice in the output

To get a 3D output of depth D , use D different filters

In practice, CNNs use many filters (~ 64 to 1024)

All together, the weights are **4** dimensional:
(output depth, input depth, kernel height, kernel width)

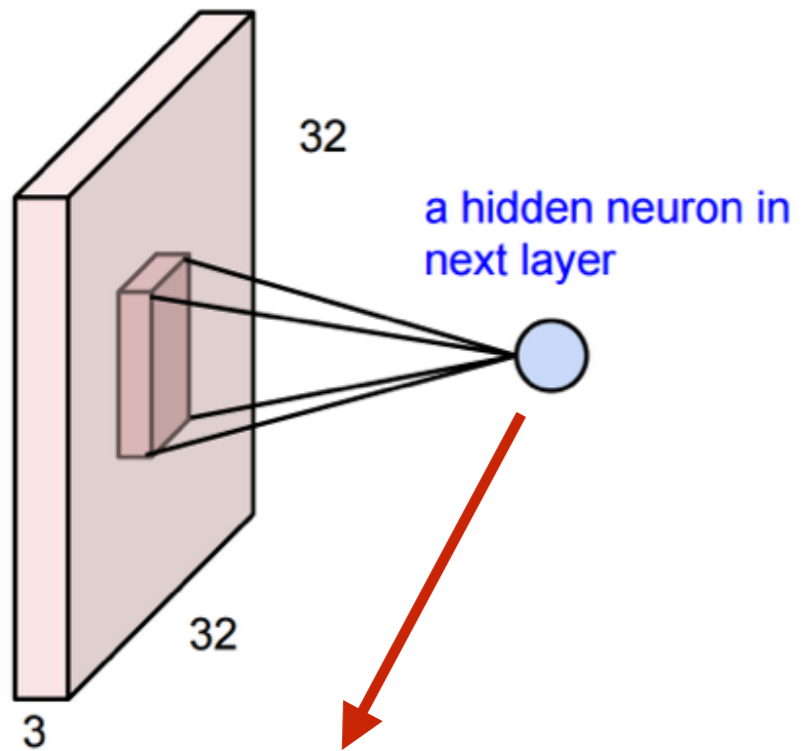
3D Activations



Let's code this up in NumPy

```
out[n, 0, r, c] =
```


3D Activations

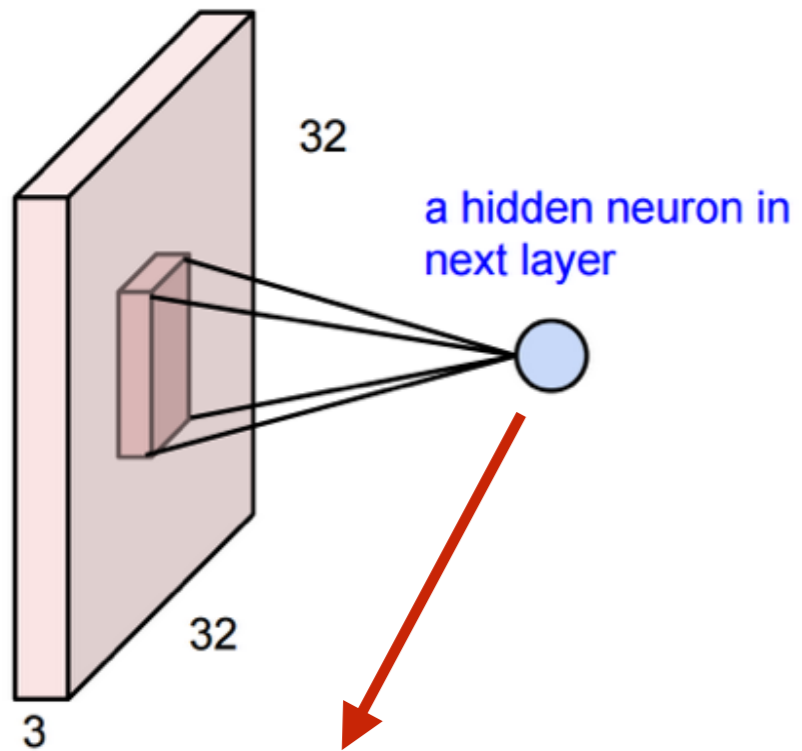


Let's code this up in NumPy

```
out[n, 0, r, c] =
```

n^{th} example

3D Activations



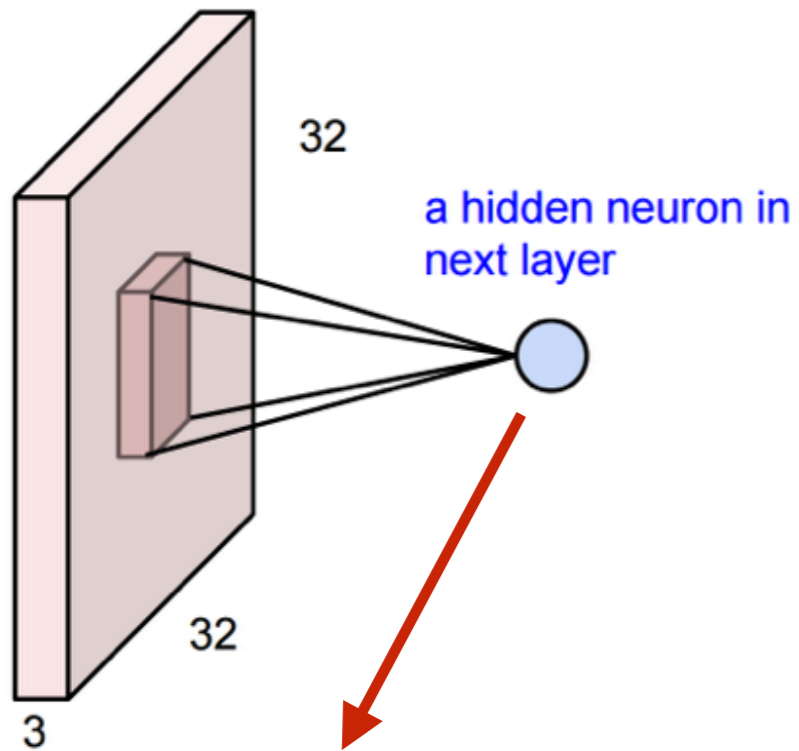
Let's code this up in NumPy

```
out[n, 0, r, c] =
```

↑
nth example

↑
first filter

3D Activations



Let's code this up in NumPy

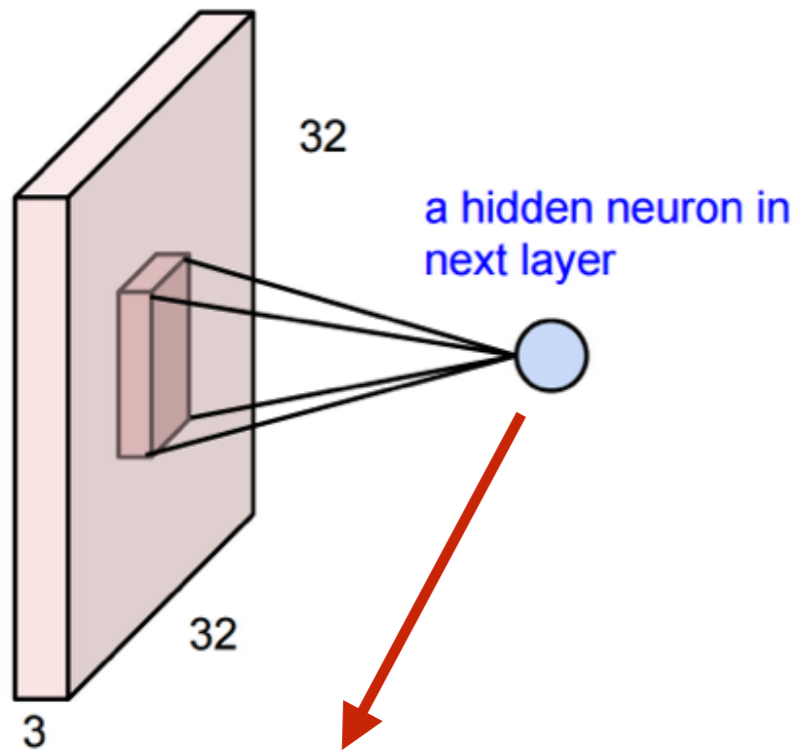
```
out[n, 0, r, c] =
```

↑
nth example

↑
first filter

↑ ↑
output position

3D Activations



Let's code this up in NumPy

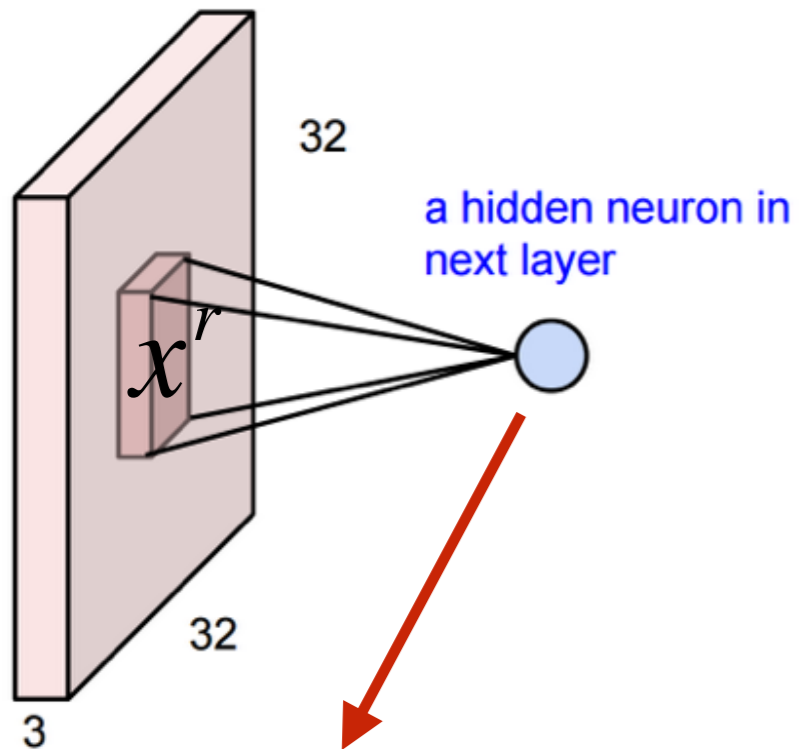
```
out[n, 0, r, c] = np.sum(
```

↑
nth example

↑
first filter

↑ ↑
output position

3D Activations



Let's code this up in NumPy

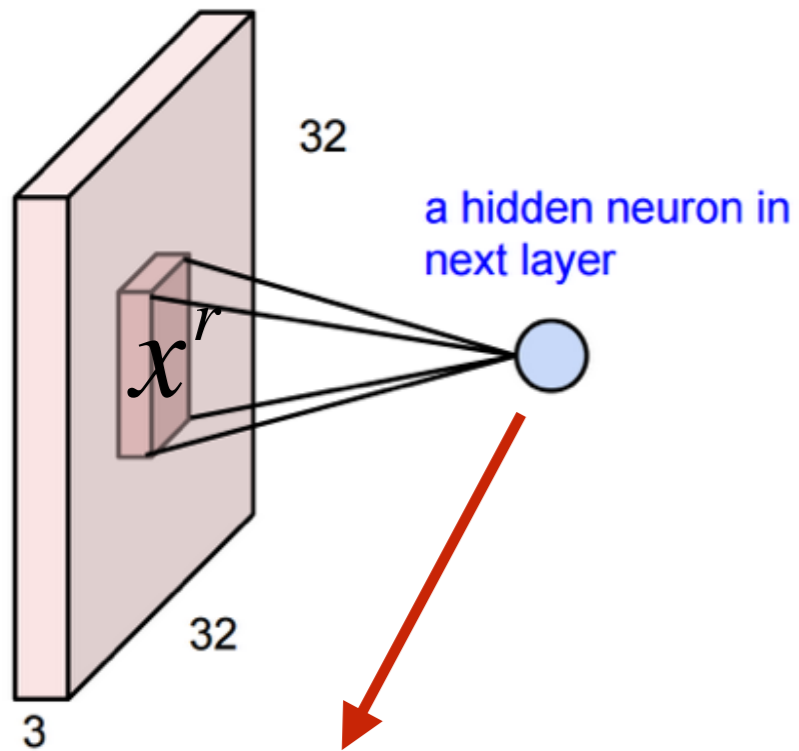
```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1])
```

↑
nth example

↑
first filter

↑ ↑
output position

3D Activations



Let's code this up in NumPy

```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1])
```

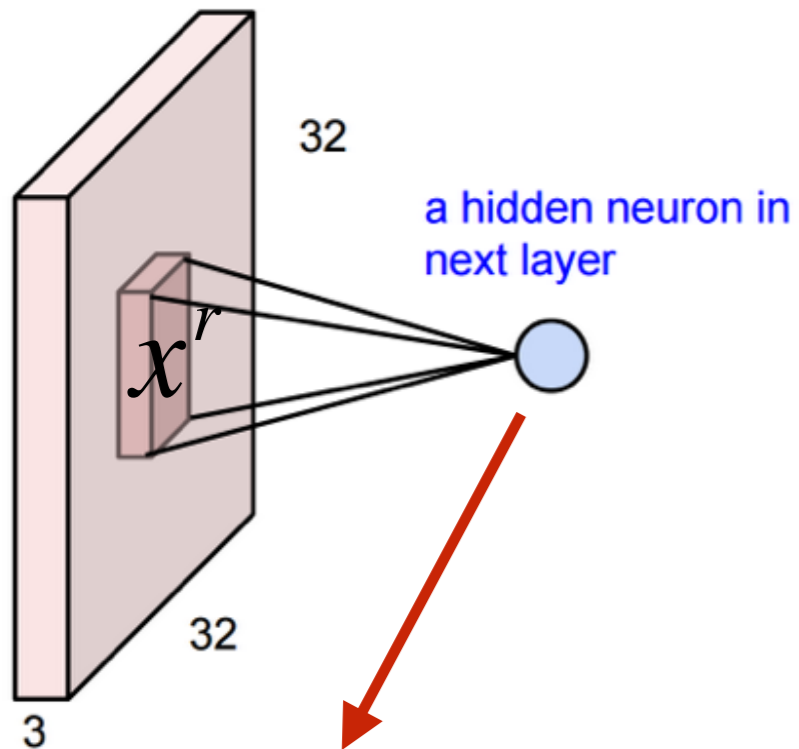
↑
nth example

↑
first filter

↑ ↑
output position

↑
nth example

3D Activations



Let's code this up in NumPy

```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1])
```

↑
nth example

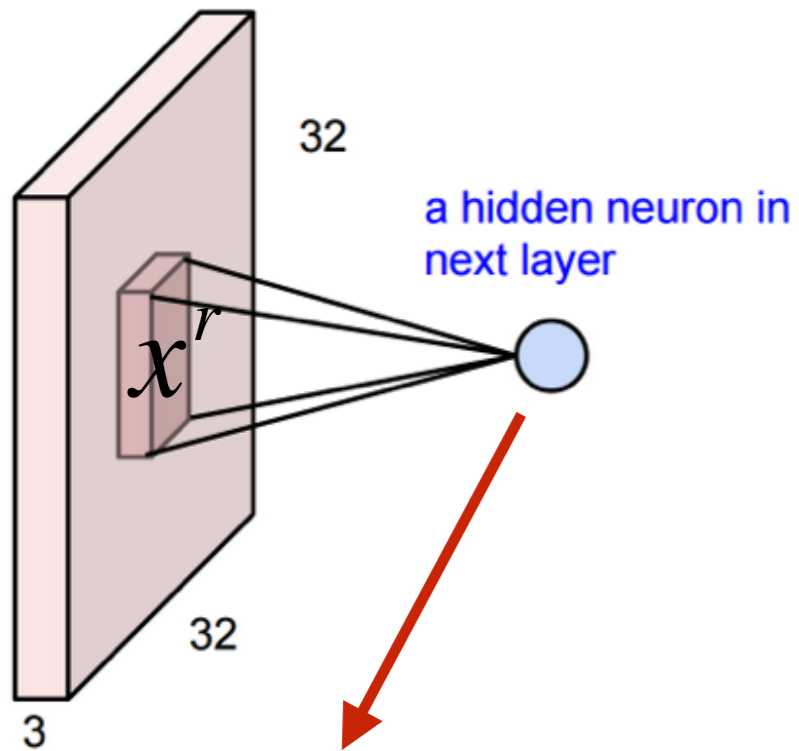
↑
first filter

↑
output position

↑
nth example

↑
all input channels

3D Activations



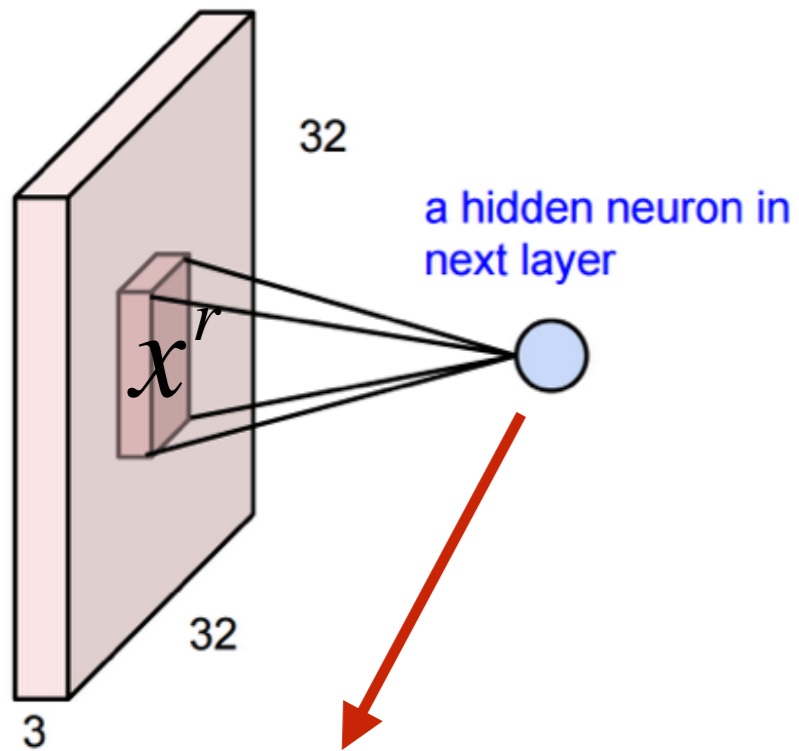
Let's code this up in NumPy

```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1])
```

n^{th} example
first filter
output position

n^{th} example
all input channels
input region

3D Activations



Let's code this up in NumPy

```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1] * W[0, :, :, :]) + b[0]
```

↑
nth example

↑
first filter

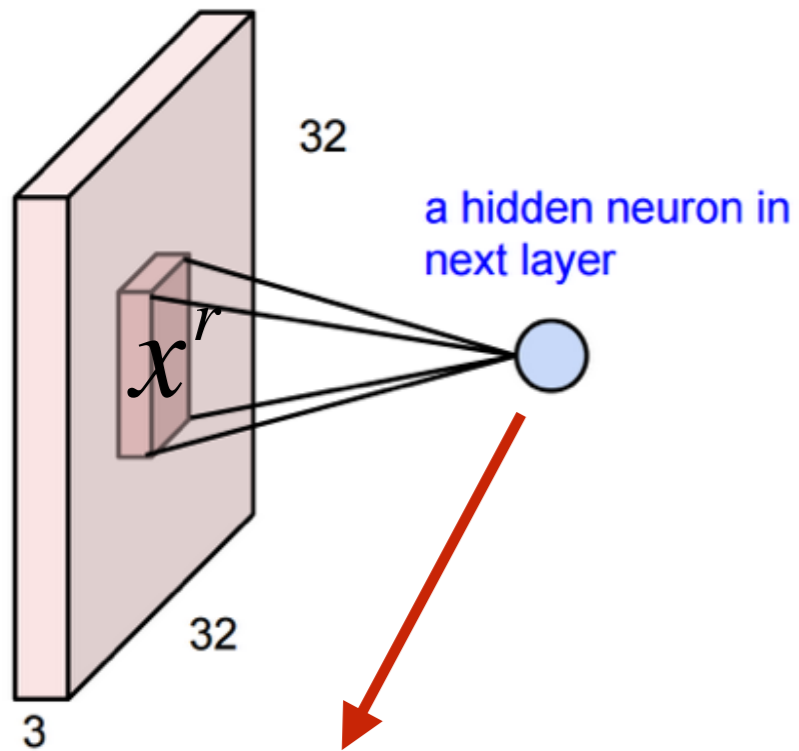
↑ ↑
output position

↑
nth example

↑
all input channels

↑ ↑
input region

3D Activations



Let's code this up in NumPy

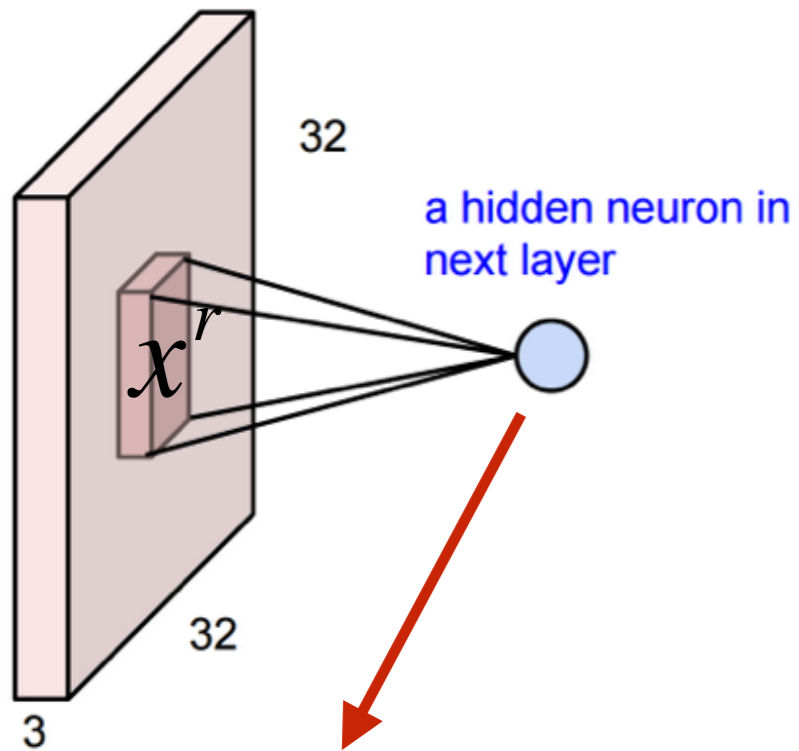
```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1] * W[0, :, :, :]) + b[0]
```

n^{th} example
first filter
output position

n^{th} example
all input channels
input region

first filter

3D Activations



Let's code this up in NumPy

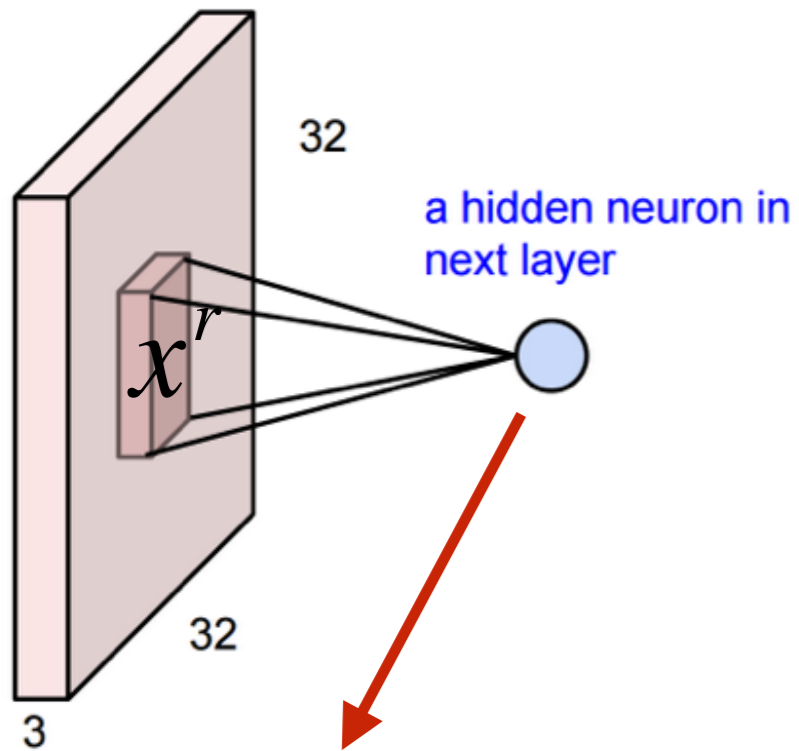
```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1] * W[0, :, :, :]) + b[0]
```

n^{th} example
first filter
output position

n^{th} example
all input channels
input region

first filter
all channels

3D Activations



Let's code this up in NumPy

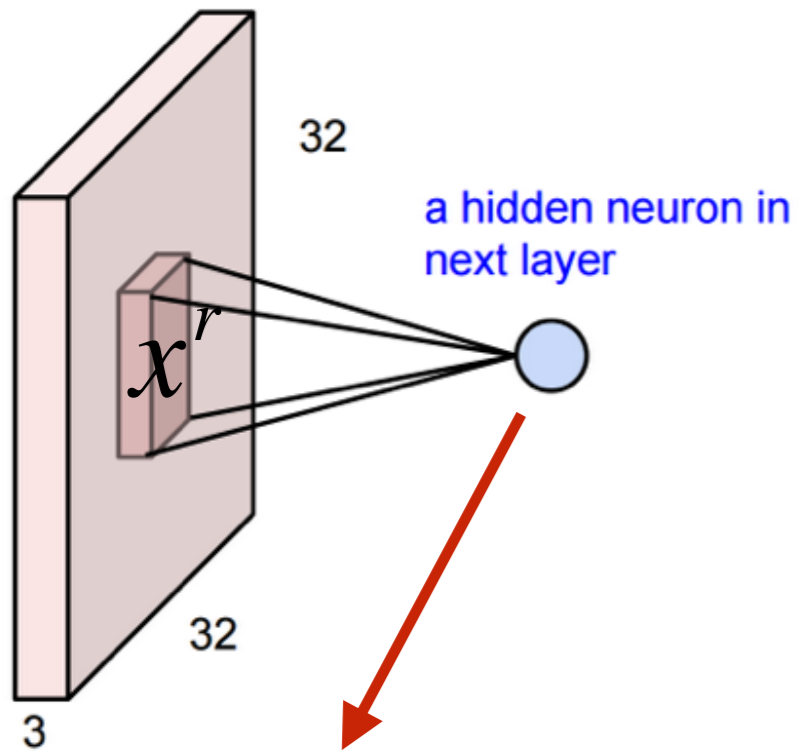
```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1] * W[0, :, :, :]) + b[0]
```

n^{th} example
first filter
output position

n^{th} example
all input channels
input region

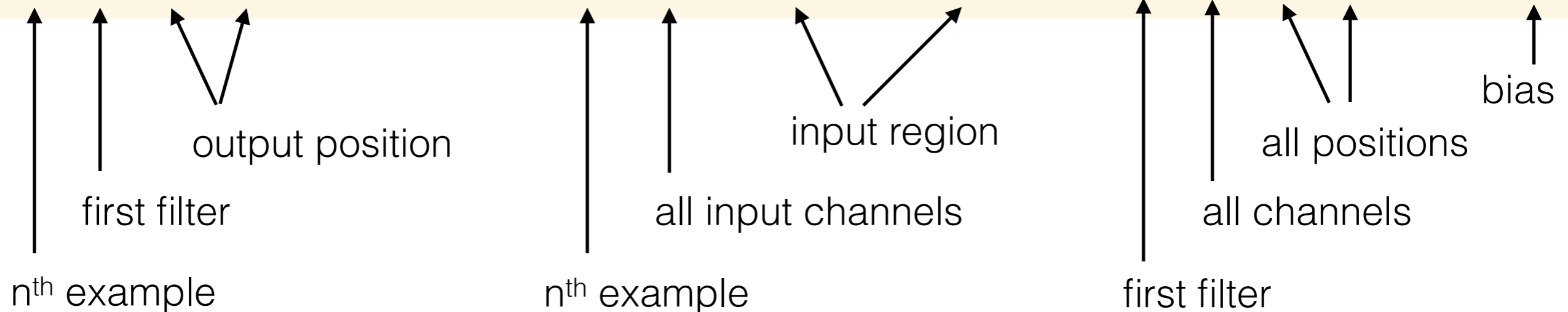
first filter
all channels
all positions

3D Activations



Let's code this up in NumPy

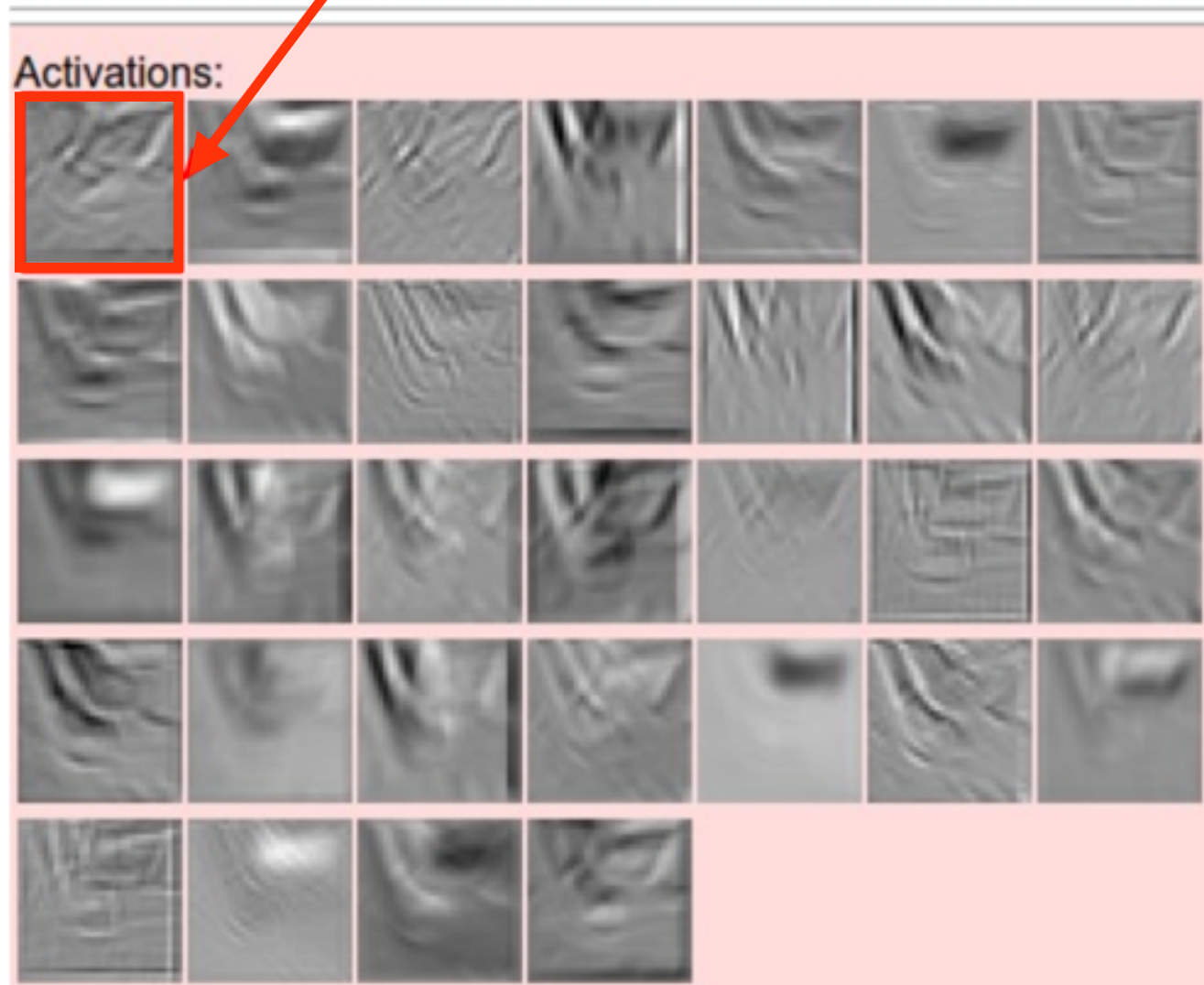
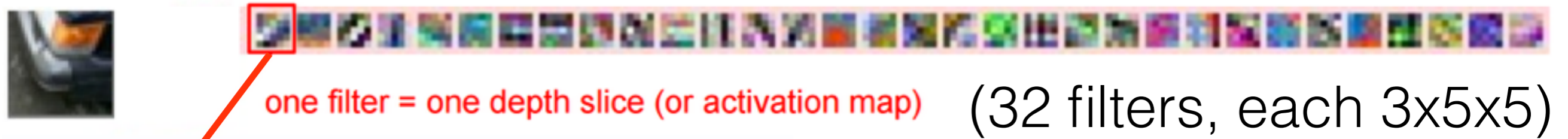
```
out[n, 0, r, c] = np.sum(X[n, :, r0:r1, c0:c1] * W[0, :, :, :]) + b[0]
```



3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)



3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

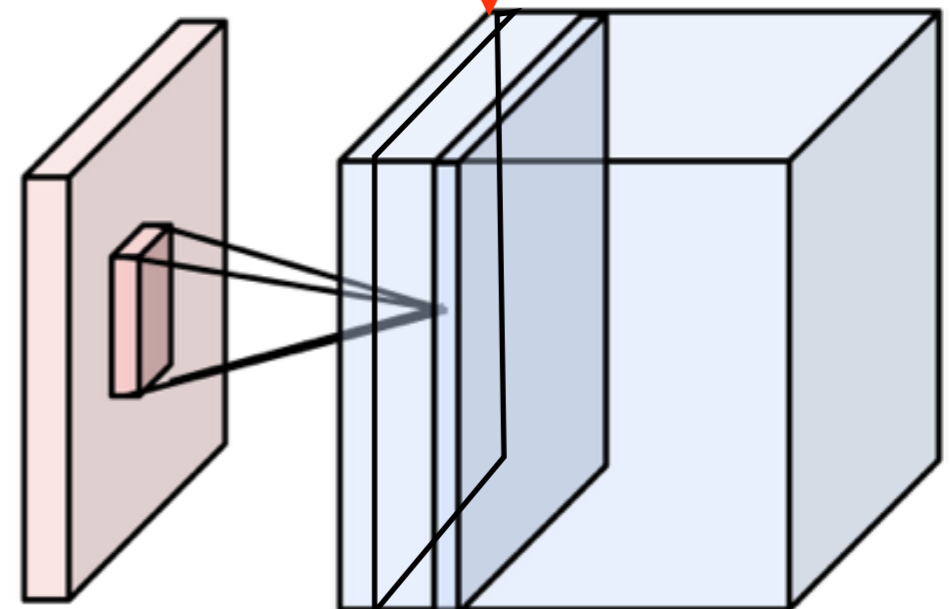
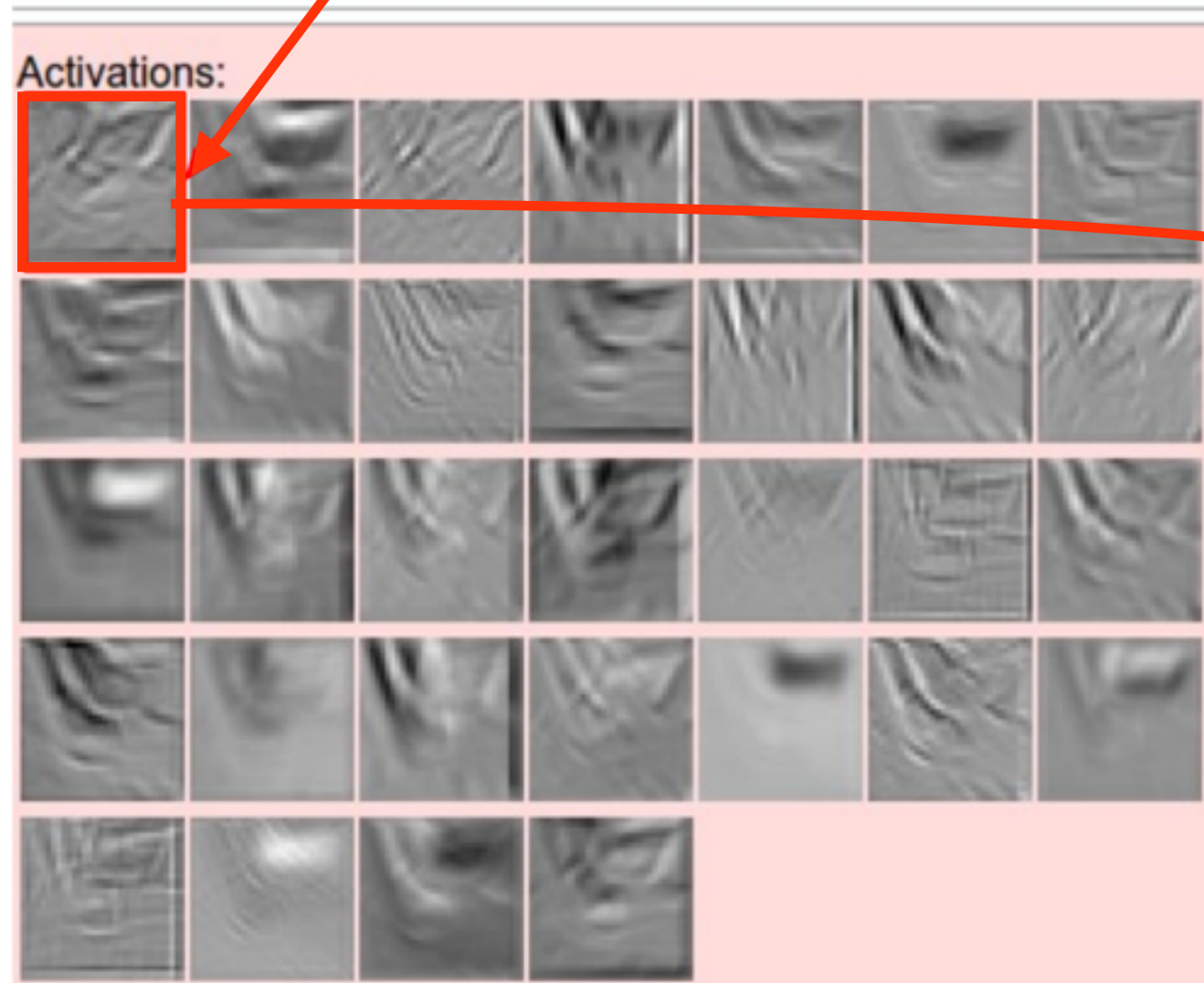
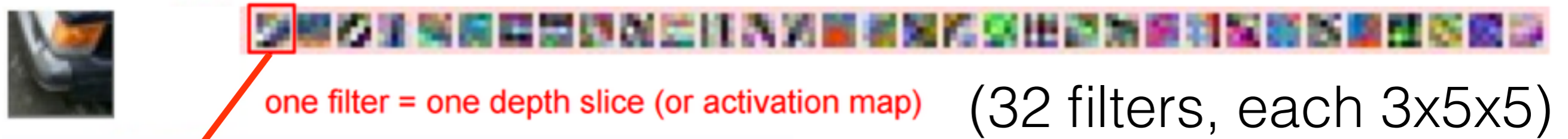


Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

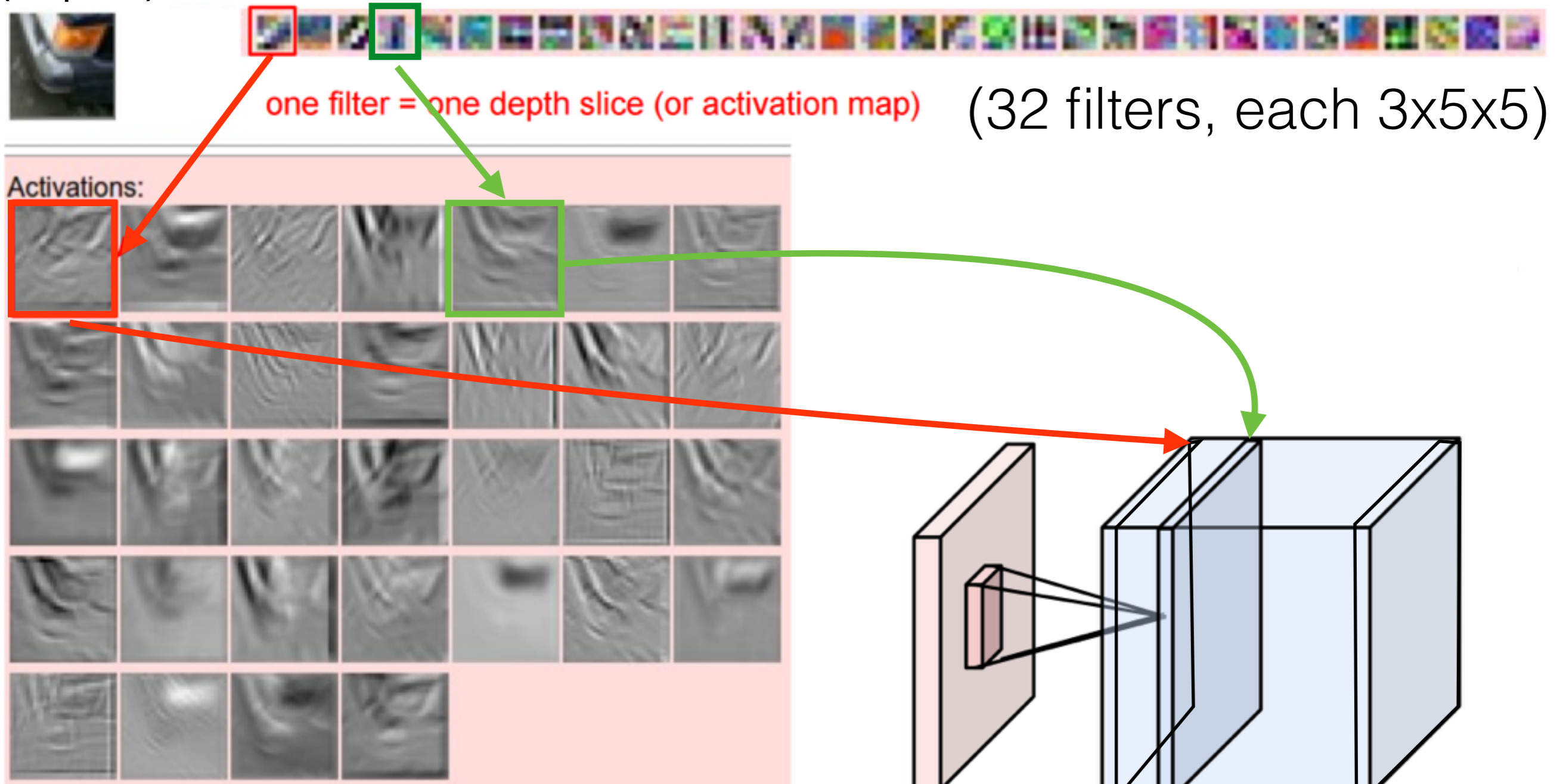
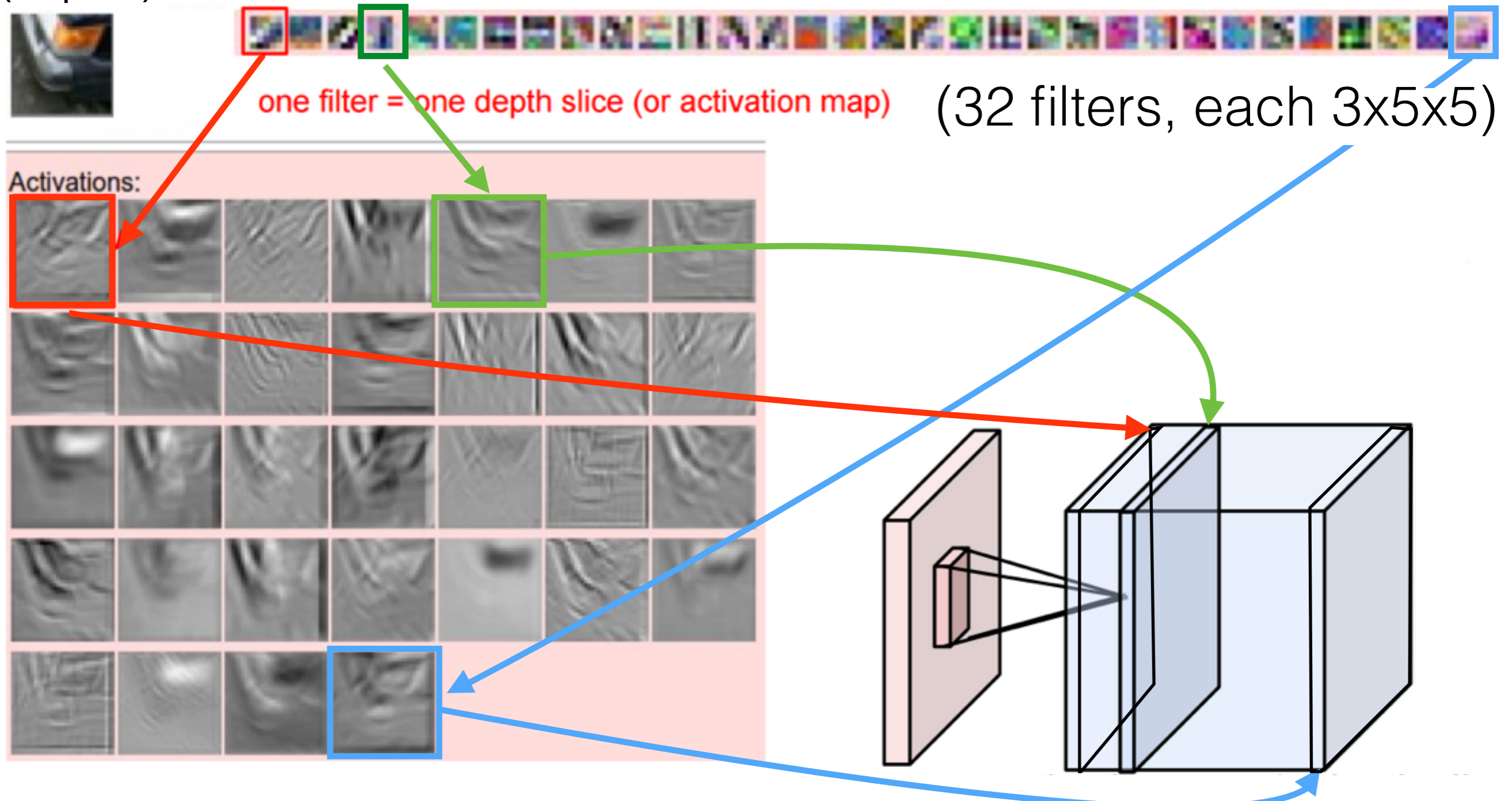


Figure: Andrej Karpathy

3D Activations

We can unravel the 3D cube and show each layer separately:

(Input)

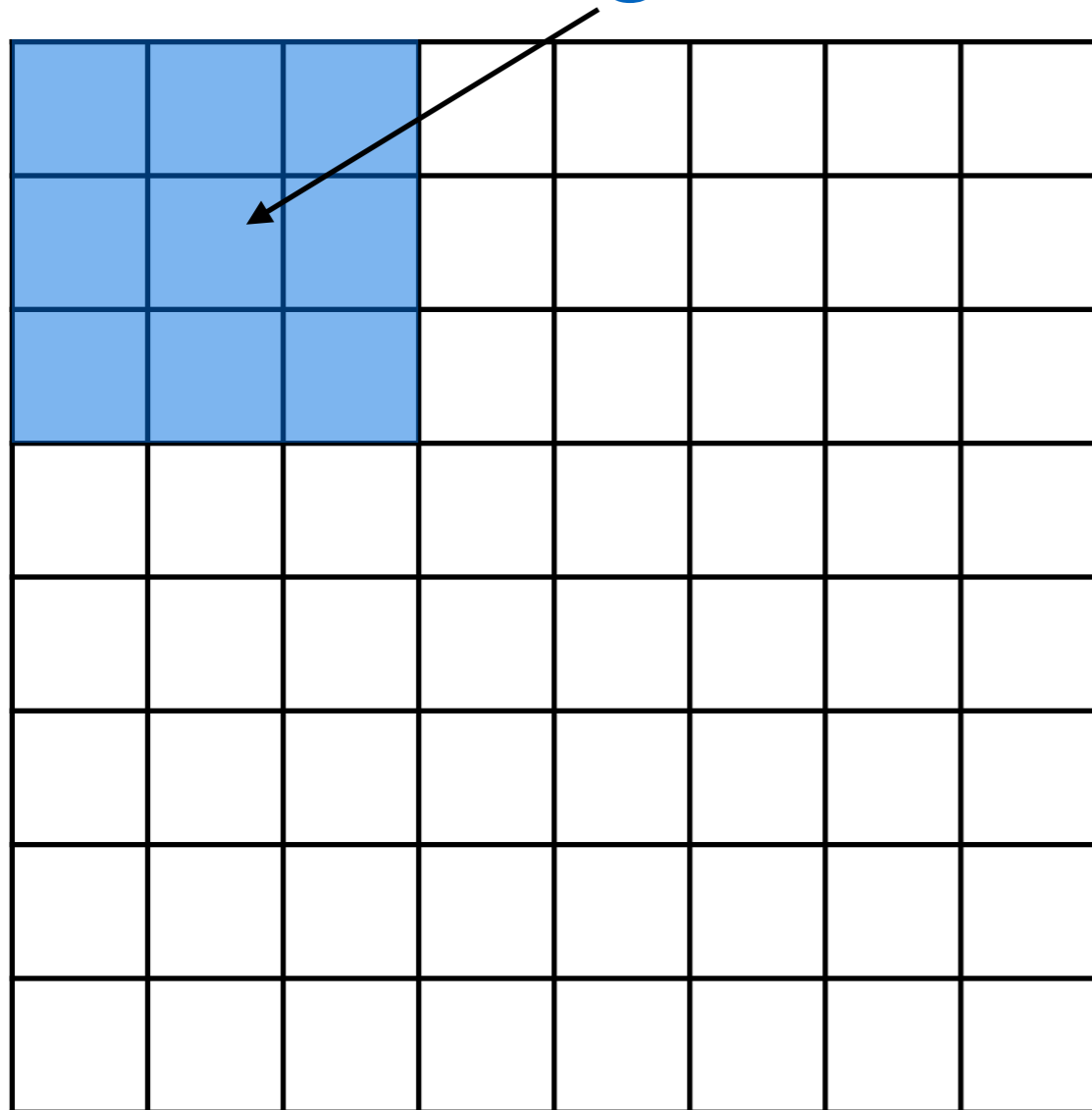


Questions?

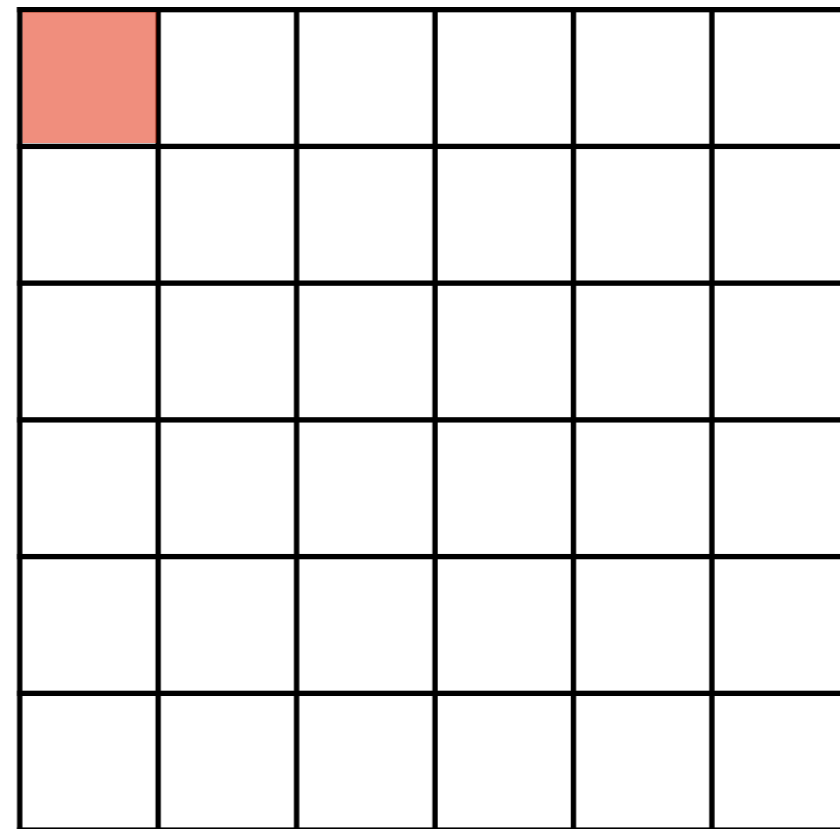
Convolution: Stride

During convolution, the weights “slide” along the input to generate each output

Weights



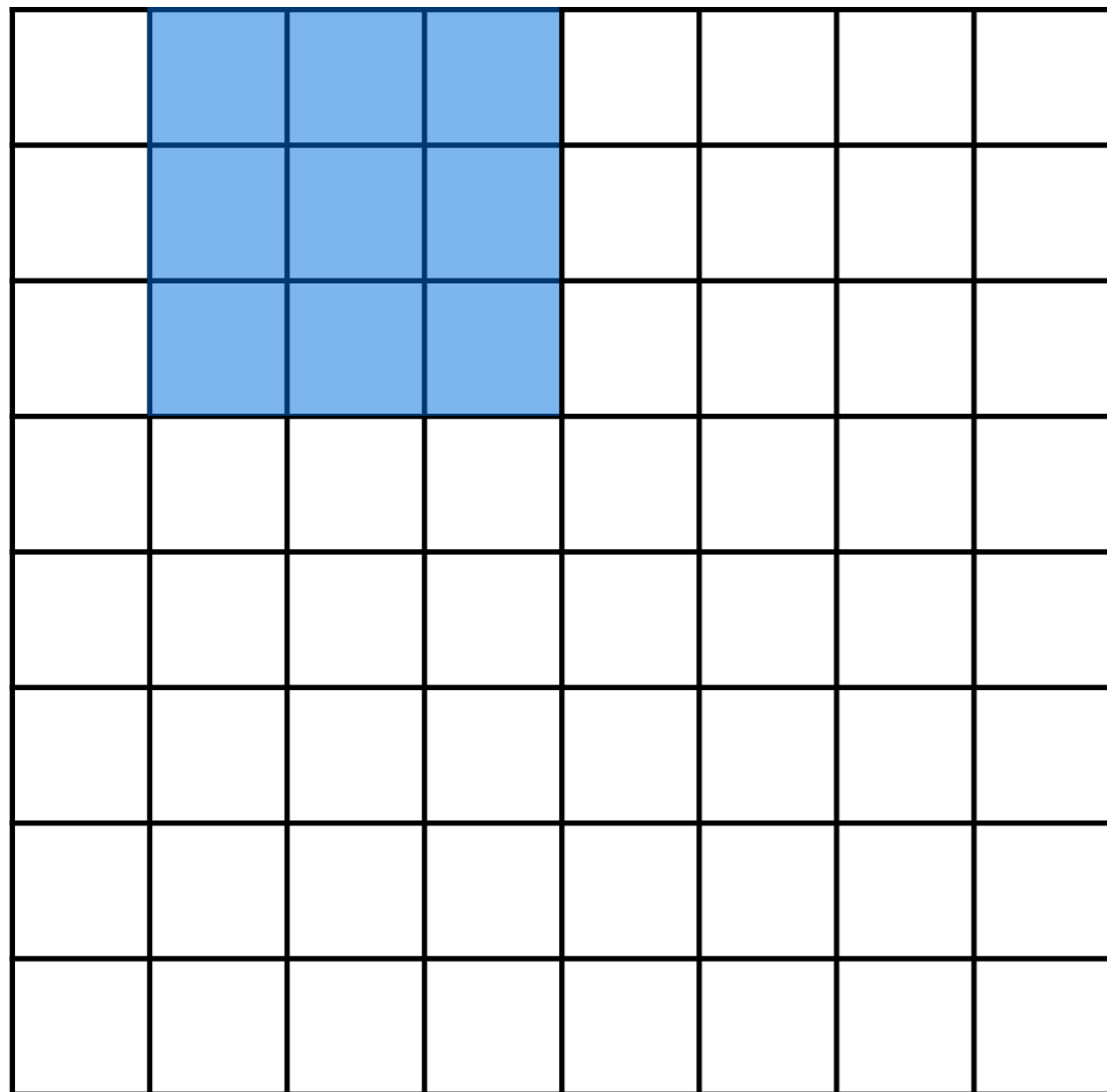
Input



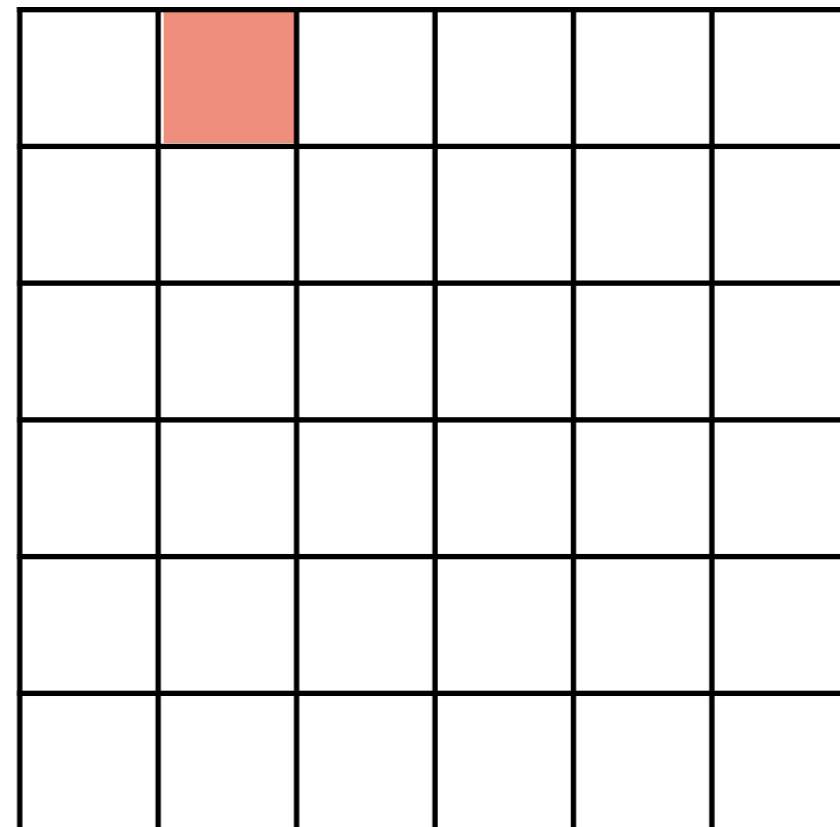
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



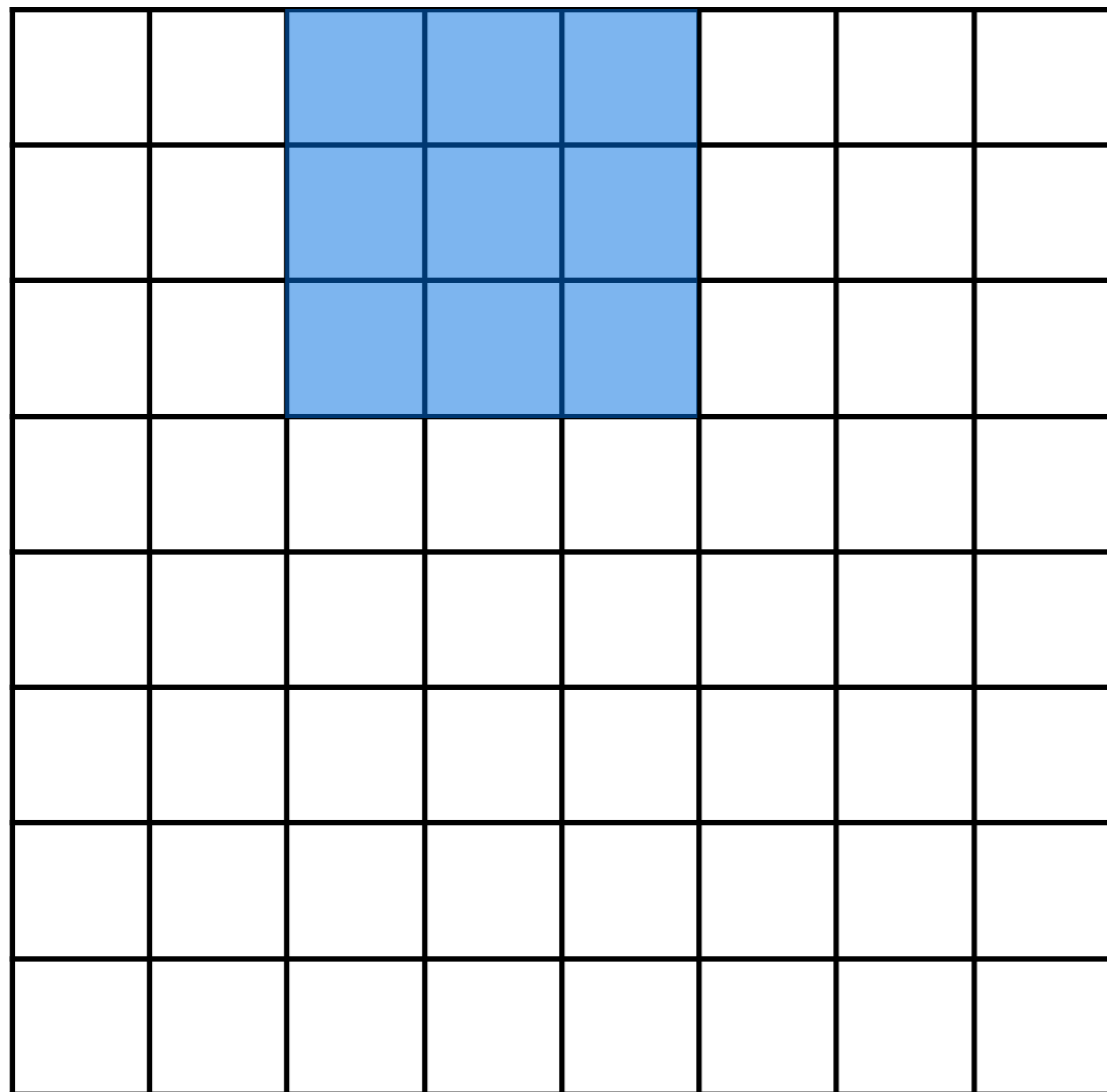
Input



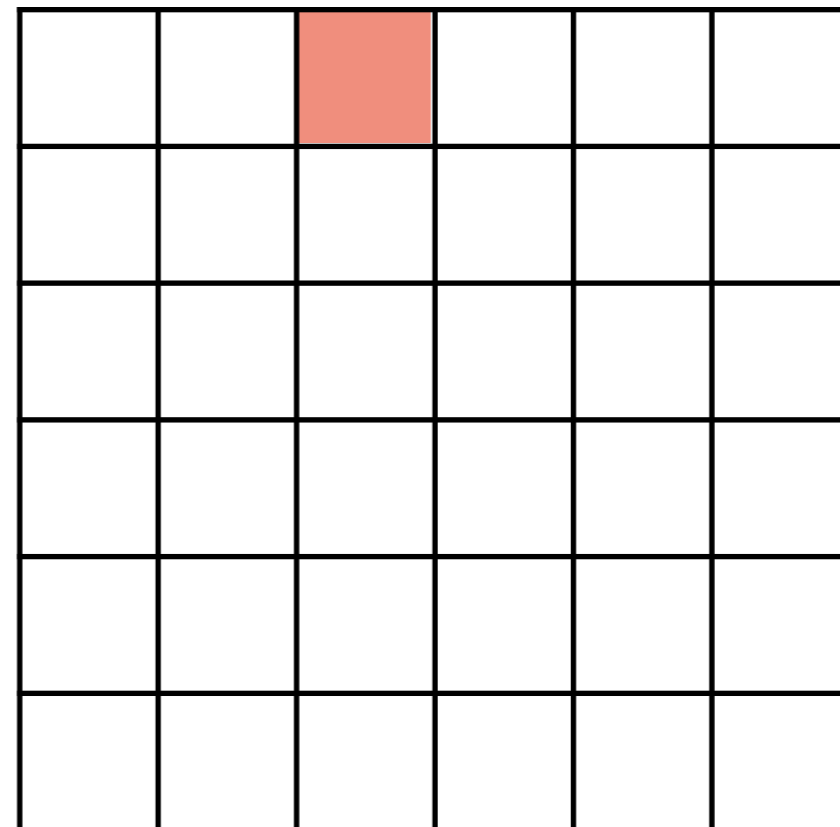
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



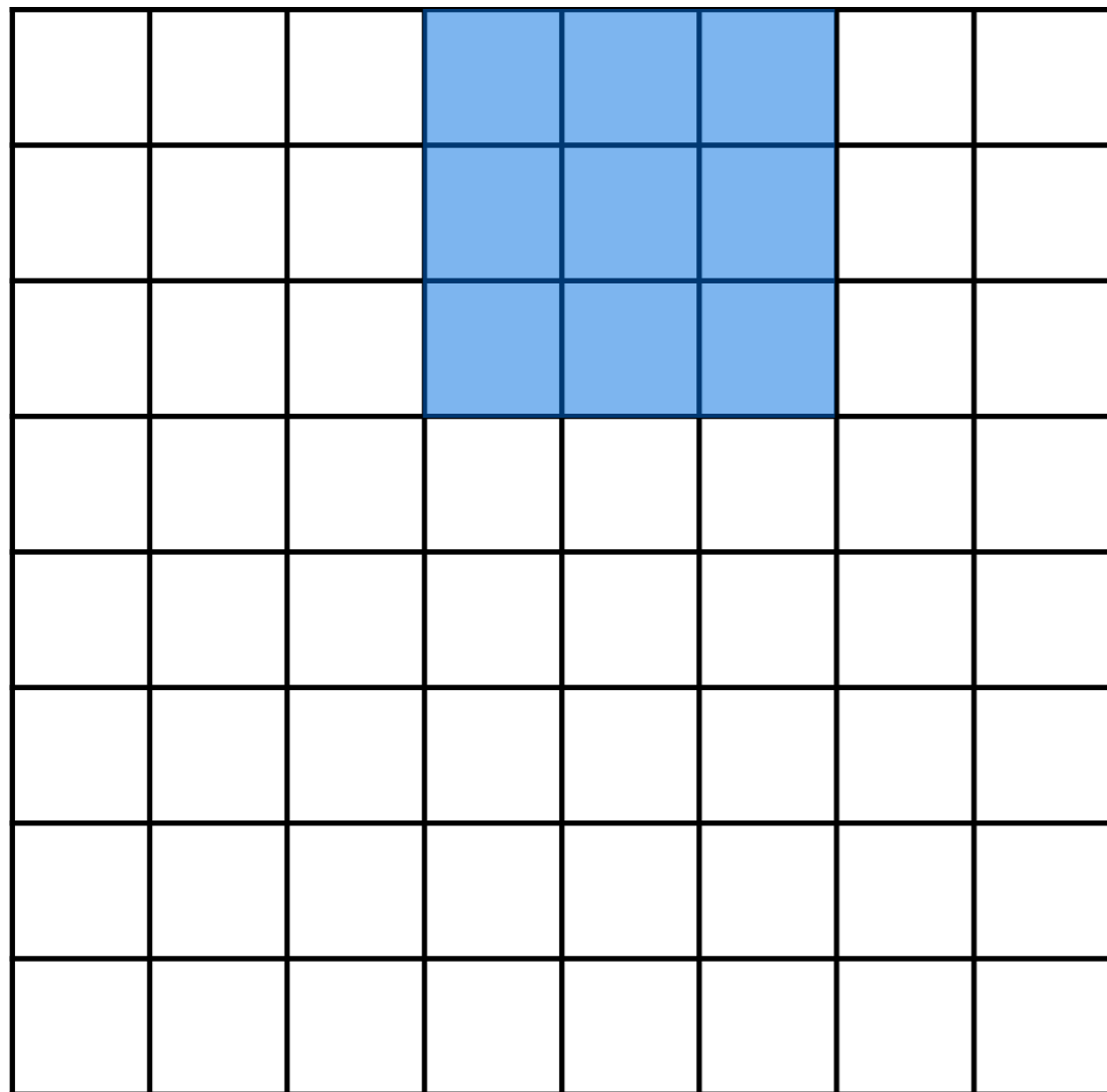
Input



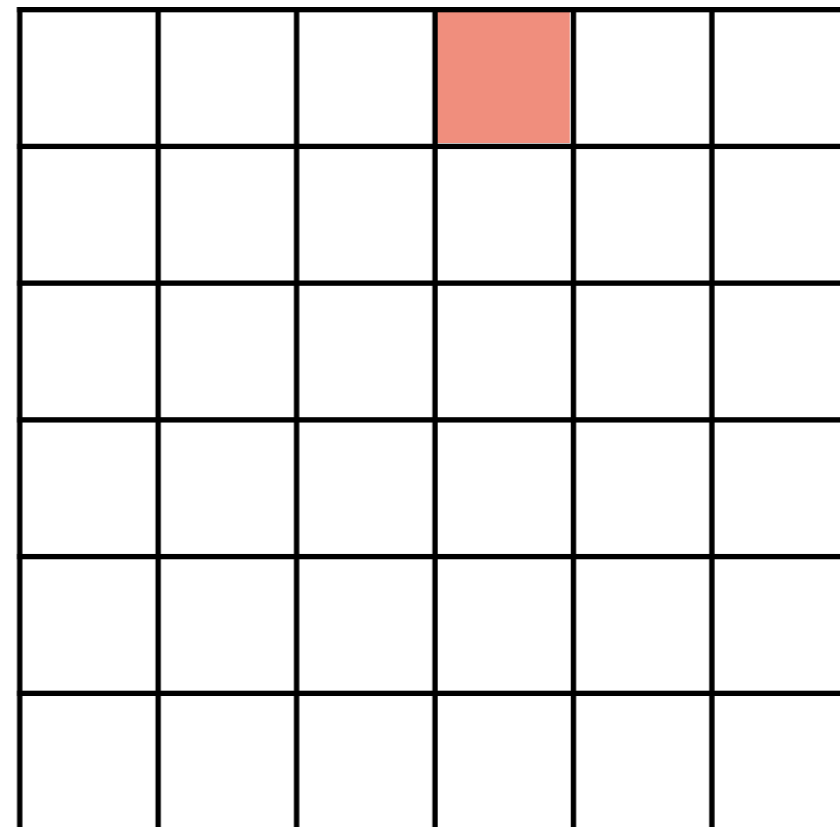
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



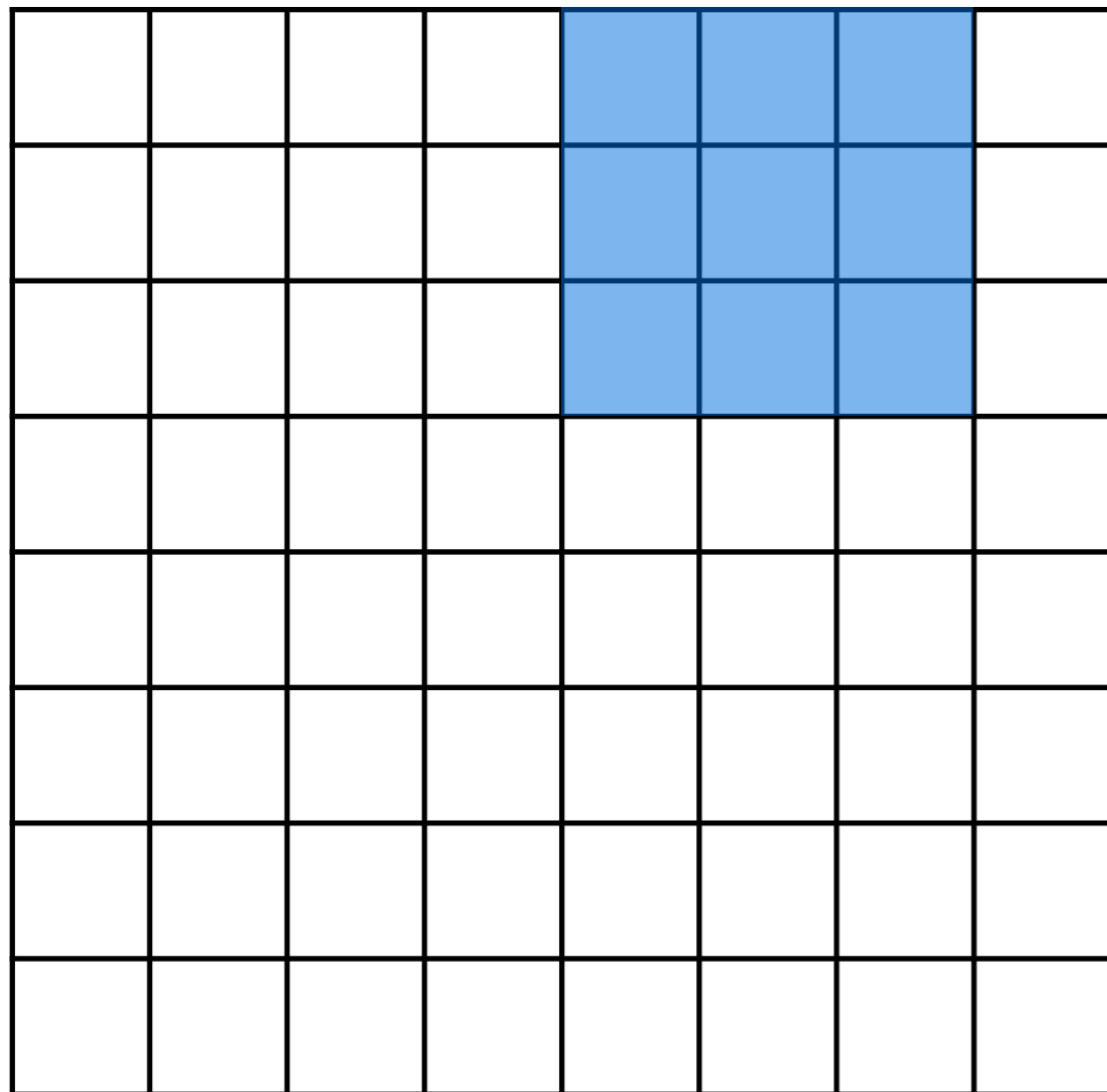
Input



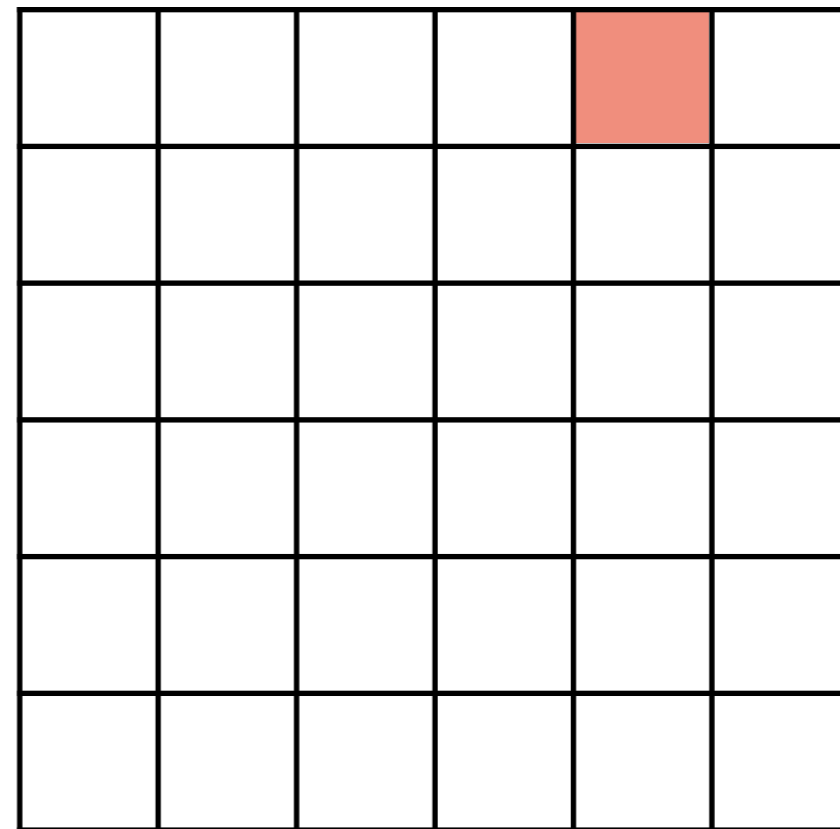
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



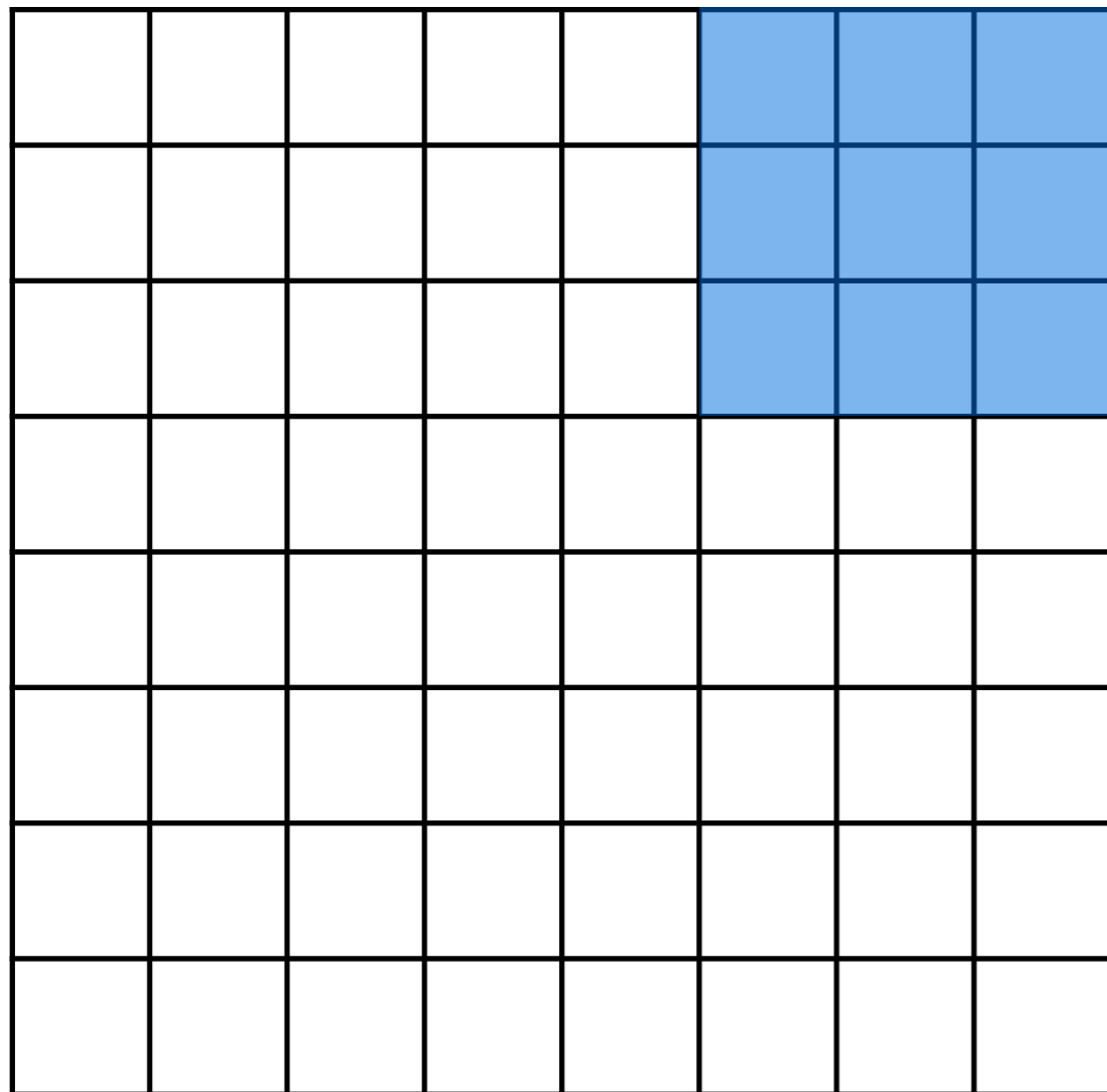
Input



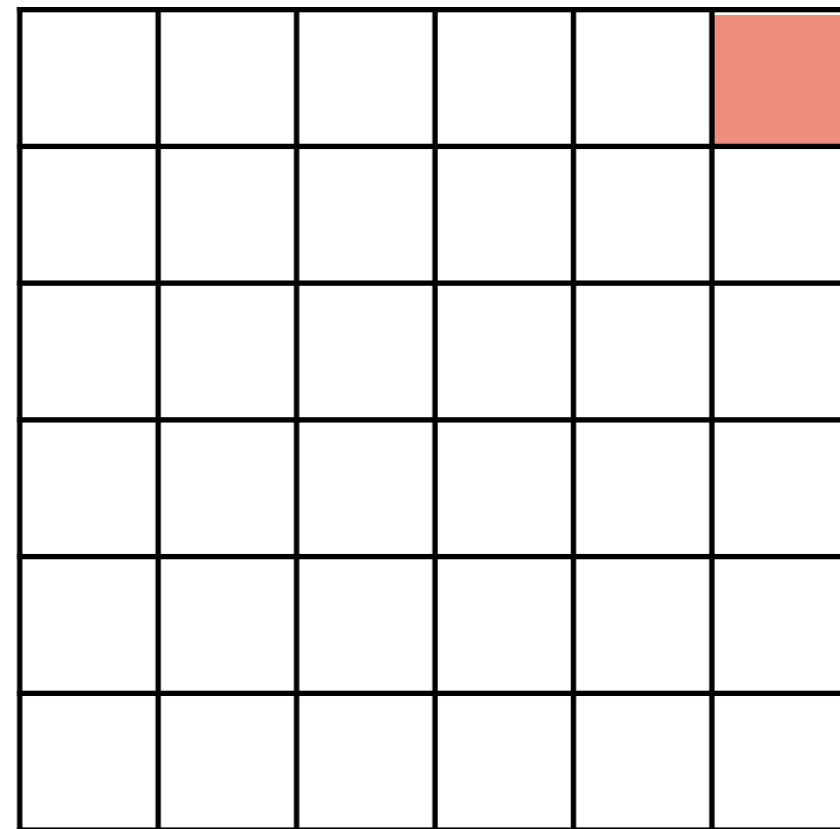
Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



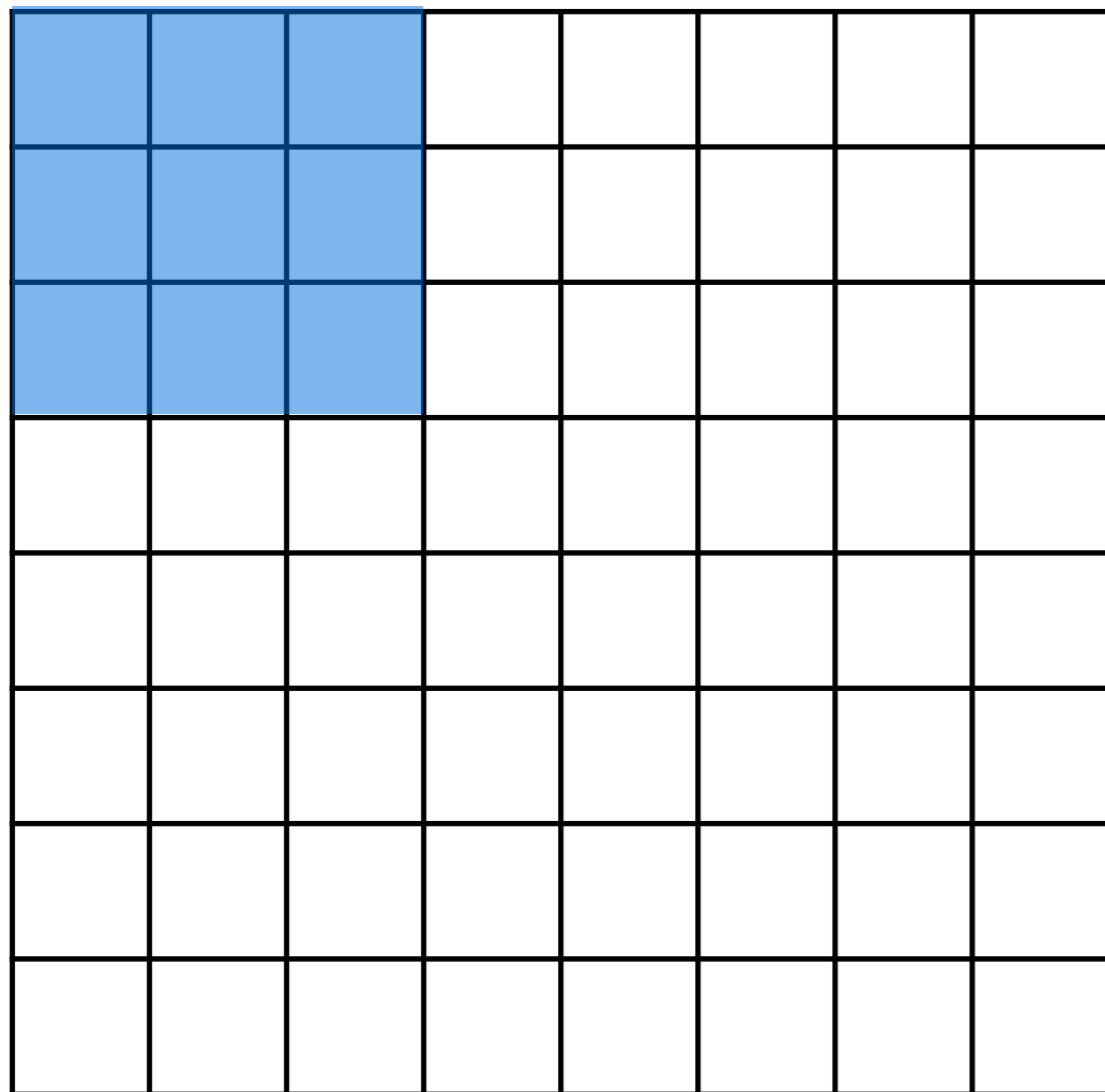
Input



Output

Convolution: Stride

During convolution, the weights “slide” along the input to generate each output



Input

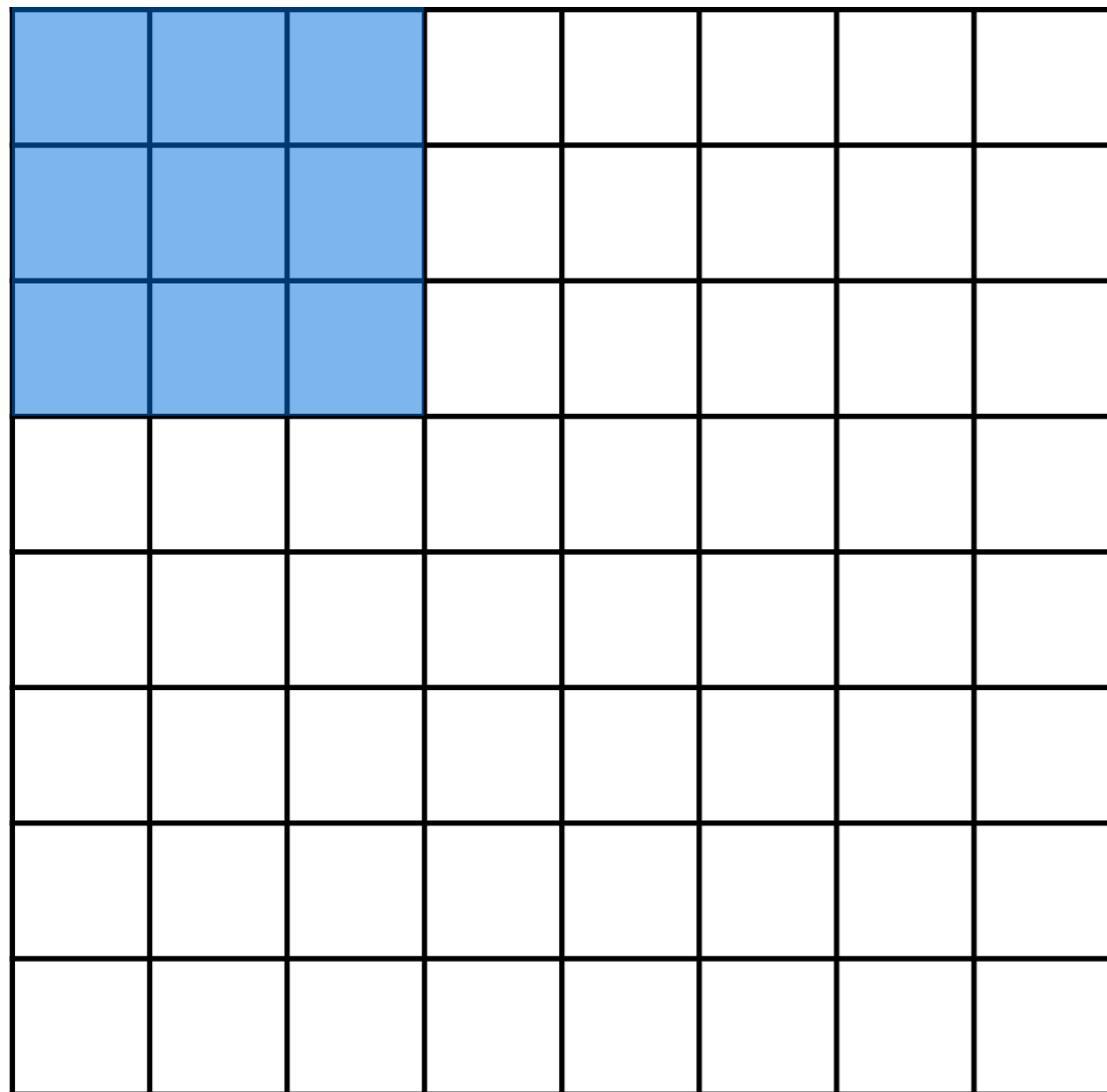
Recall that at each position, we are doing a **3D** sum:

$$h^r = \sum_{ijk} x^r_{ijk} W_{ijk} + b$$

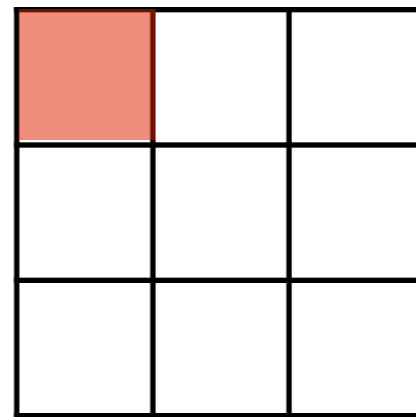
(channel, row, column)

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



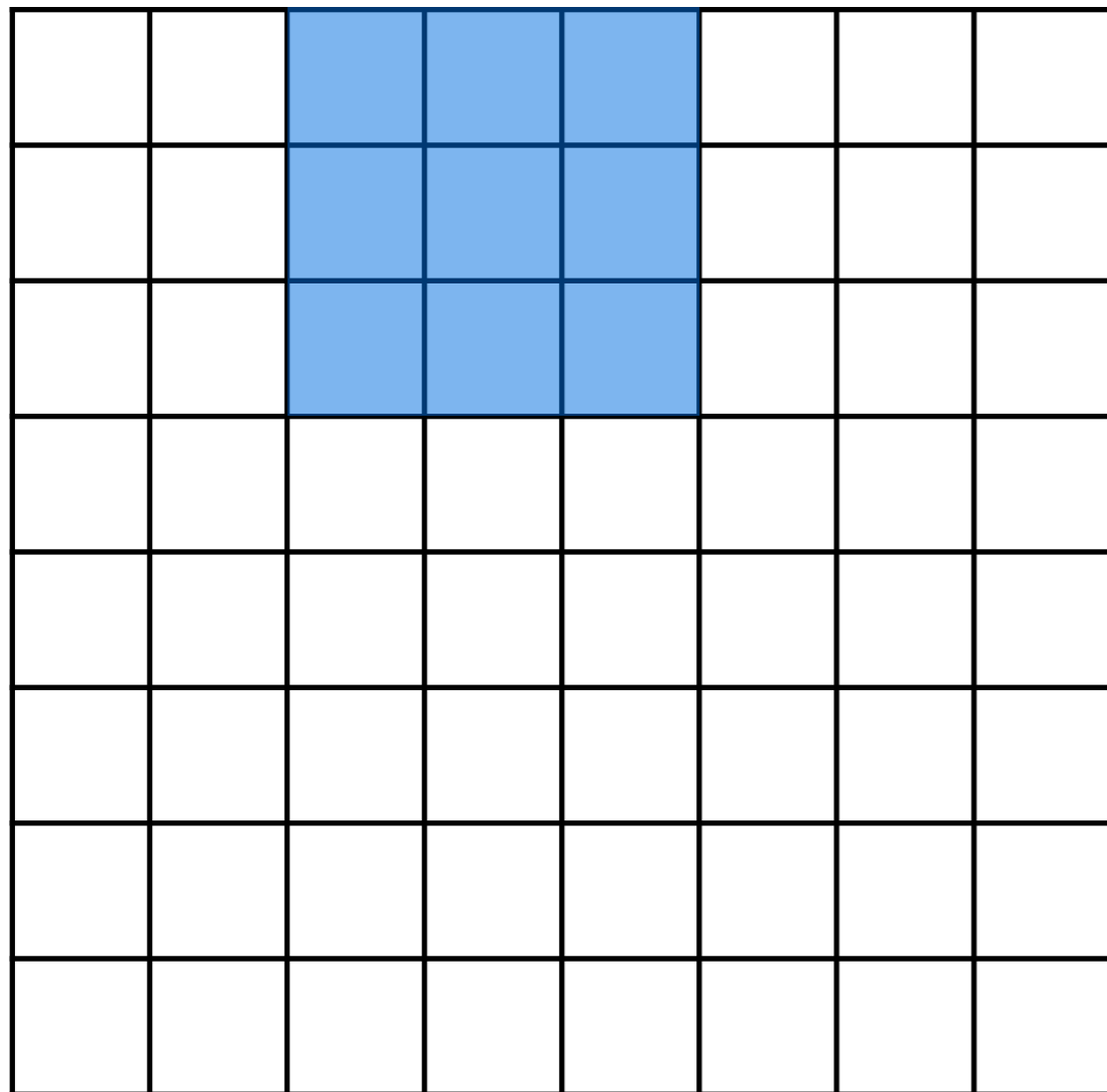
Input



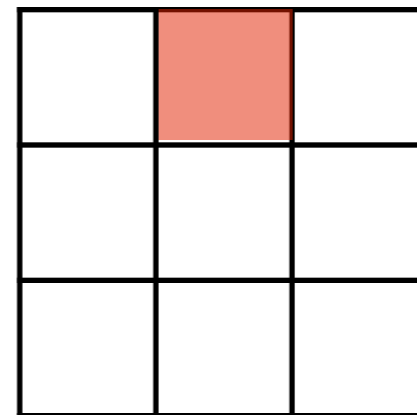
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



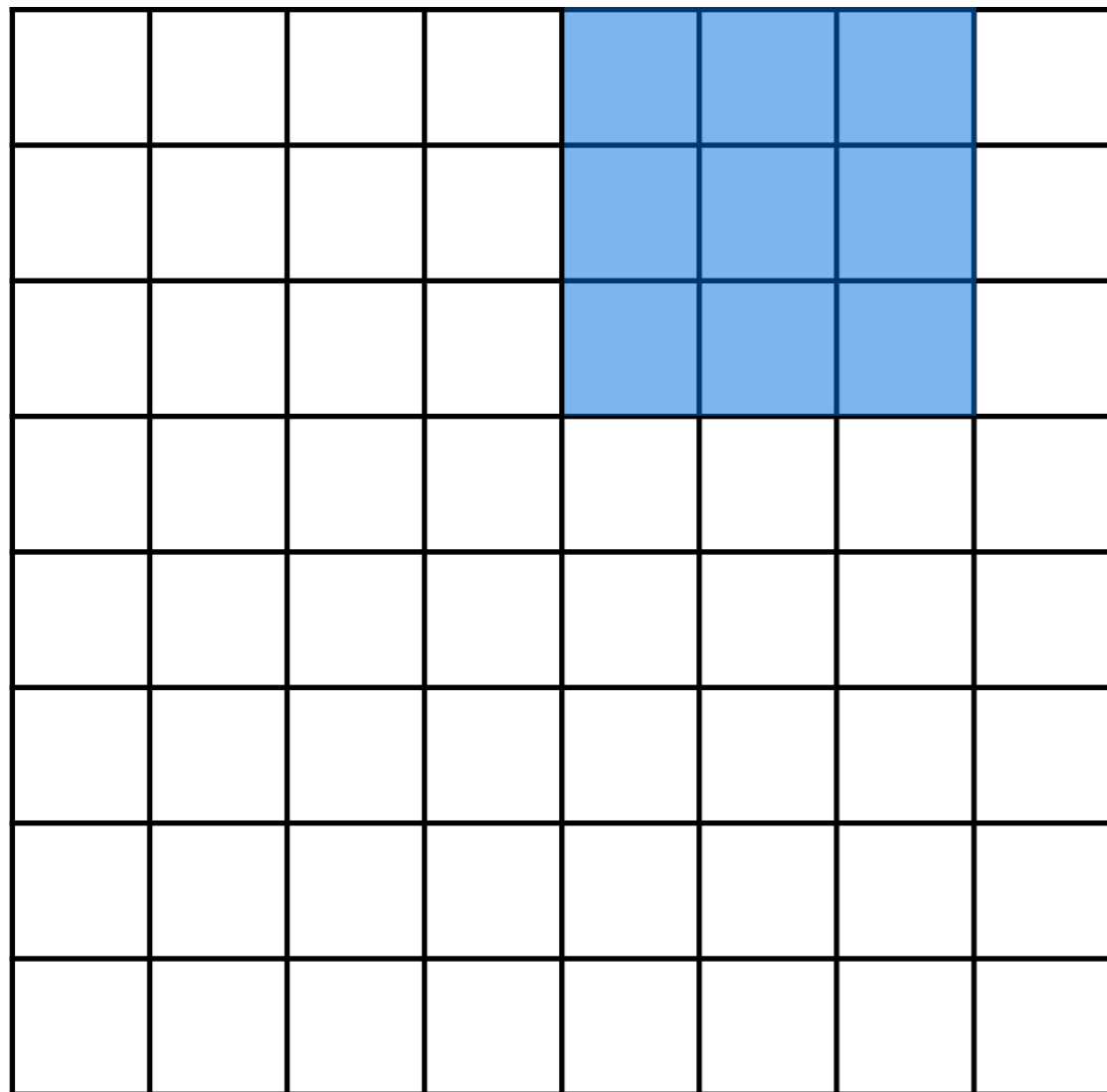
Input



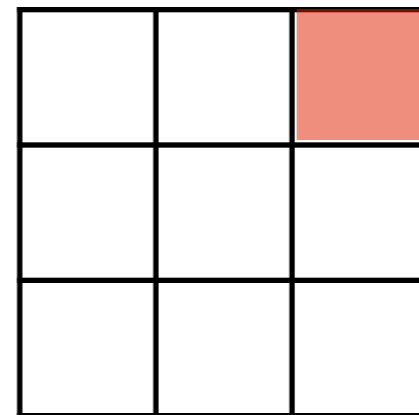
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



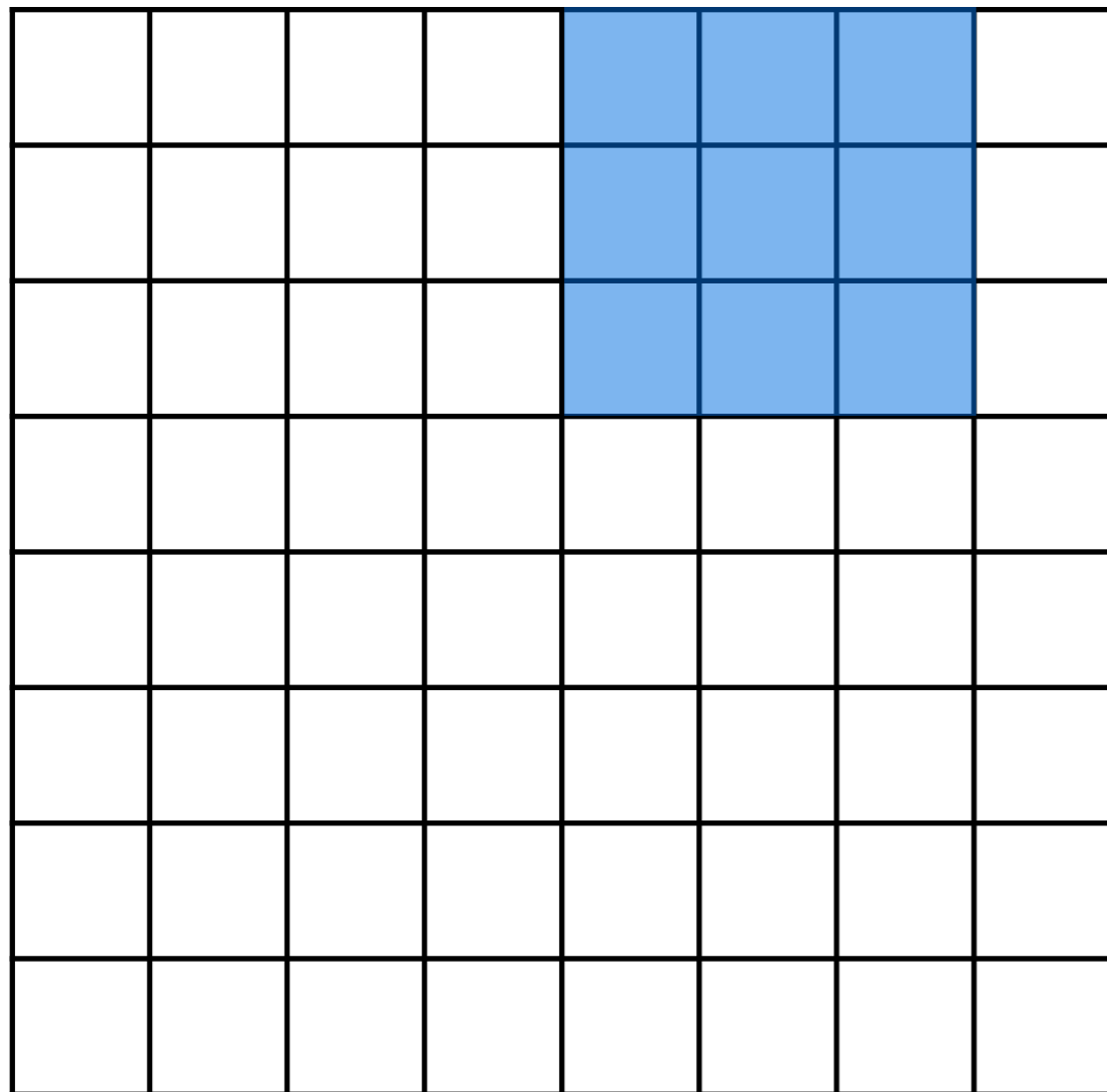
Input



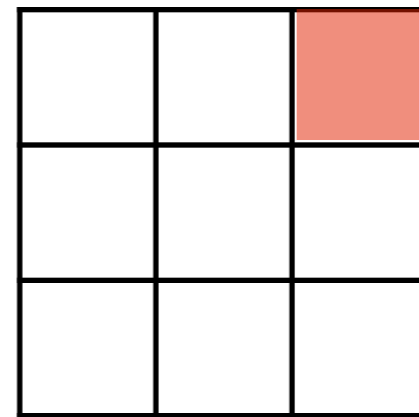
Output

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



Input

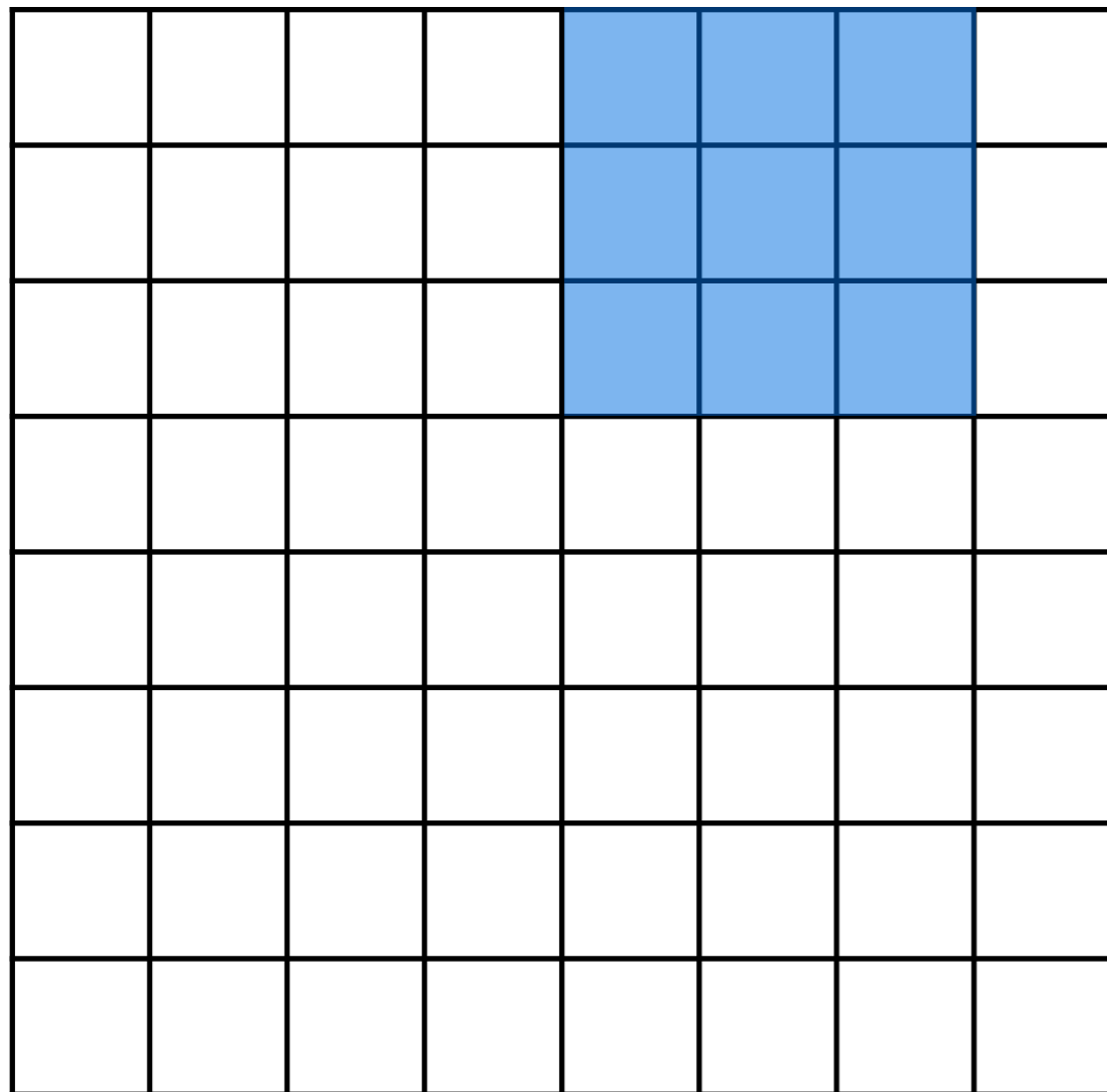


Output

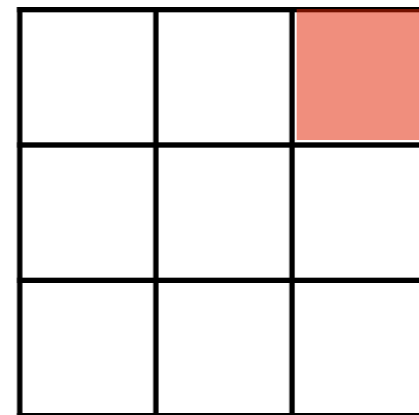
- Notice that with certain strides, we may not be able to cover all of the input

Convolution: Stride

But we can also convolve with a **stride**, e.g. stride = 2



Input



Output

- Notice that with certain strides, we may not be able to cover all of the input
- The output is also half the size of the input

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

Convolution: Padding

We can also pad the input with zeros.

Here, **pad = 1**, **stride = 2**

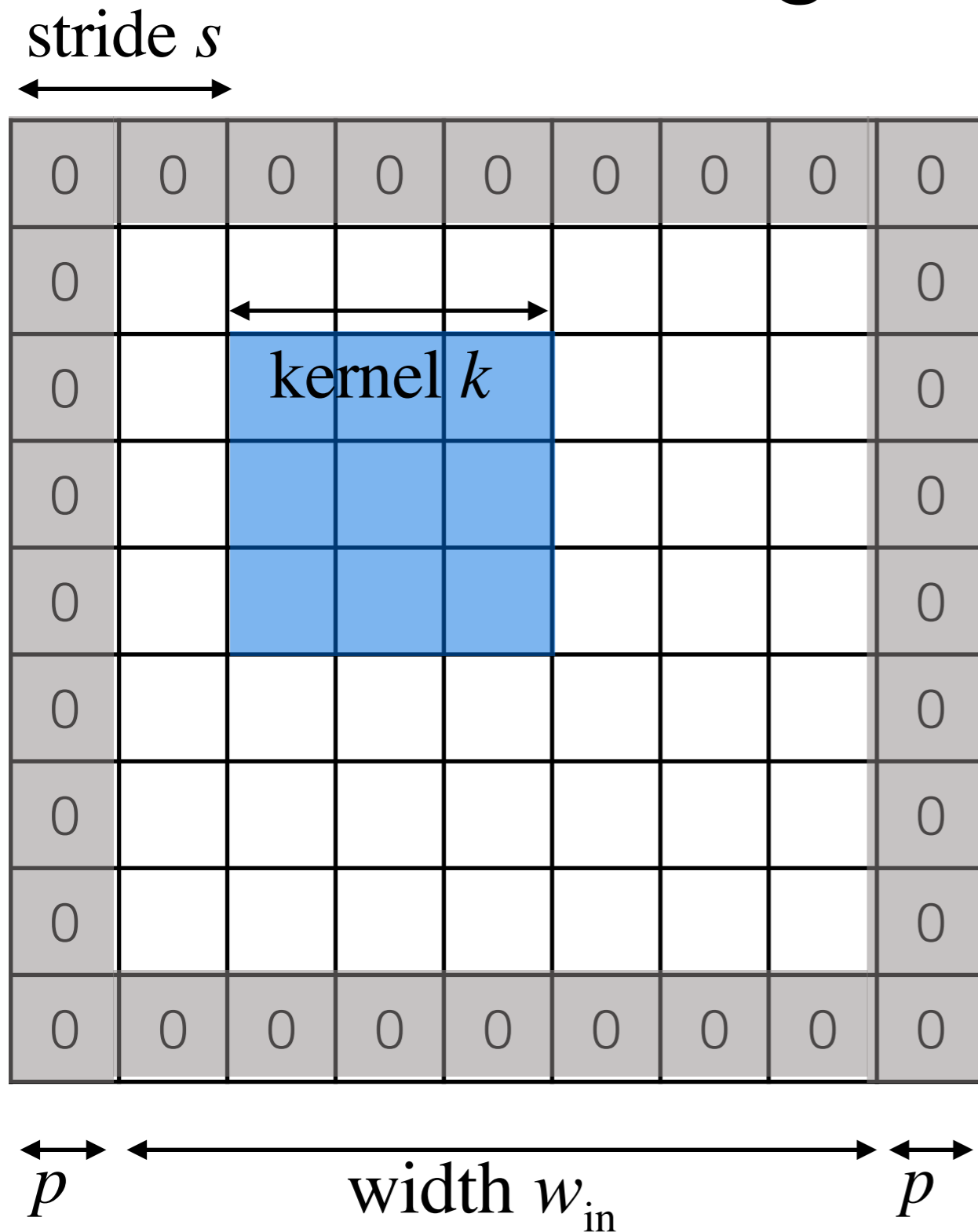
0	0	0	0	0	0	0	0	0
0						0	0	0
0						0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input

Output

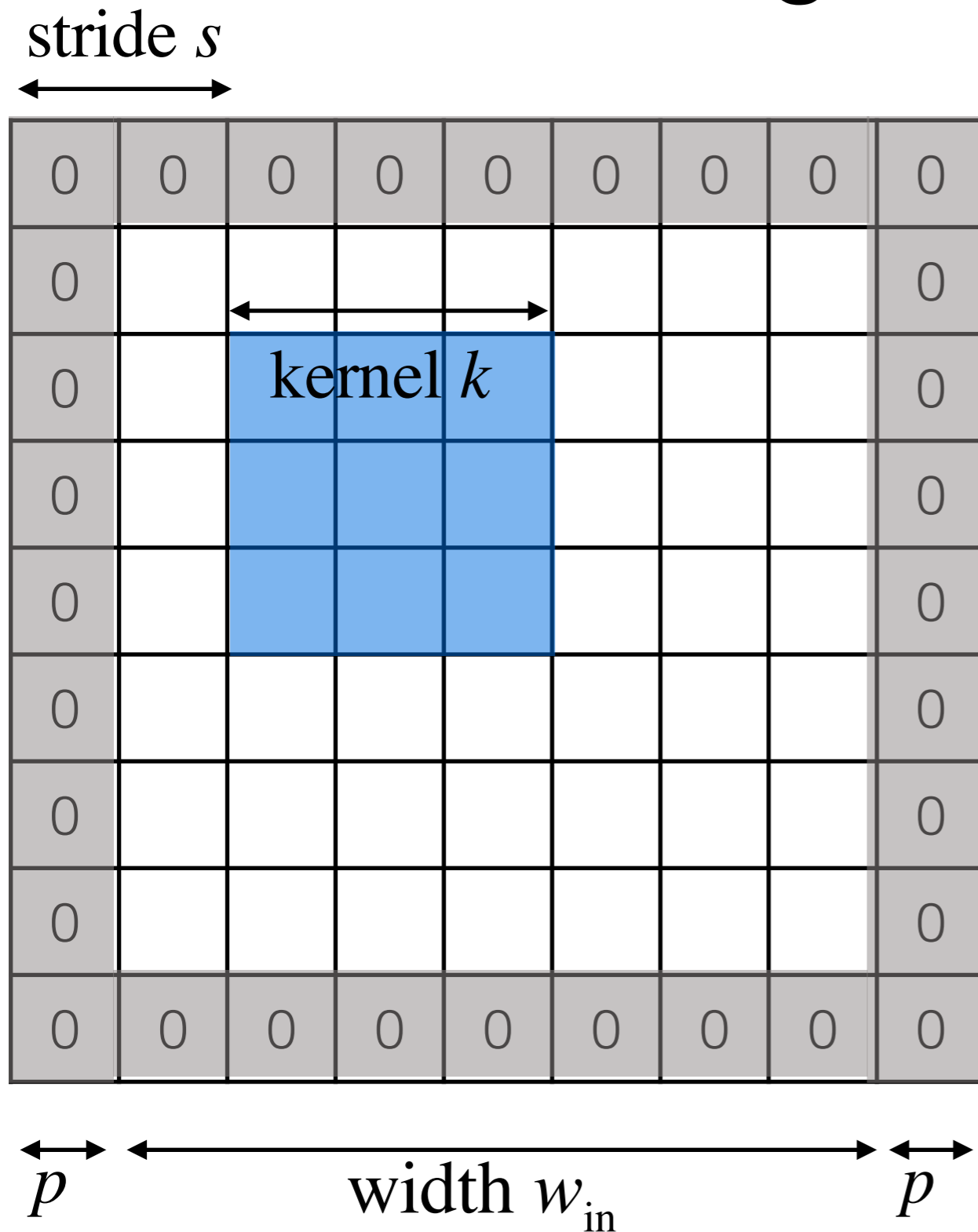
Convolution:

How big is the output?



Convolution:

How big is the output?

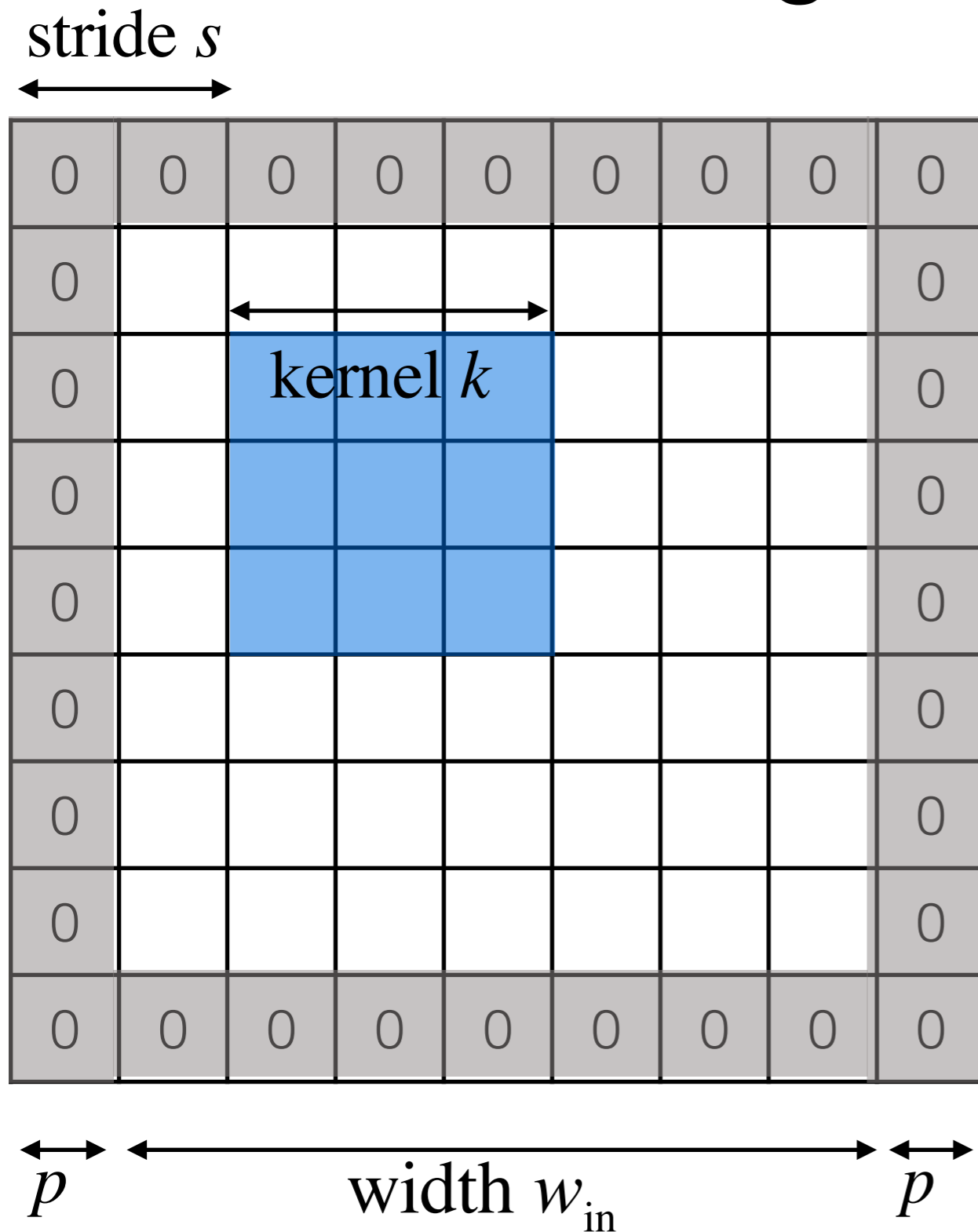


In general, the output has size:

$$w_{\text{out}} = \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1$$

Convolution:

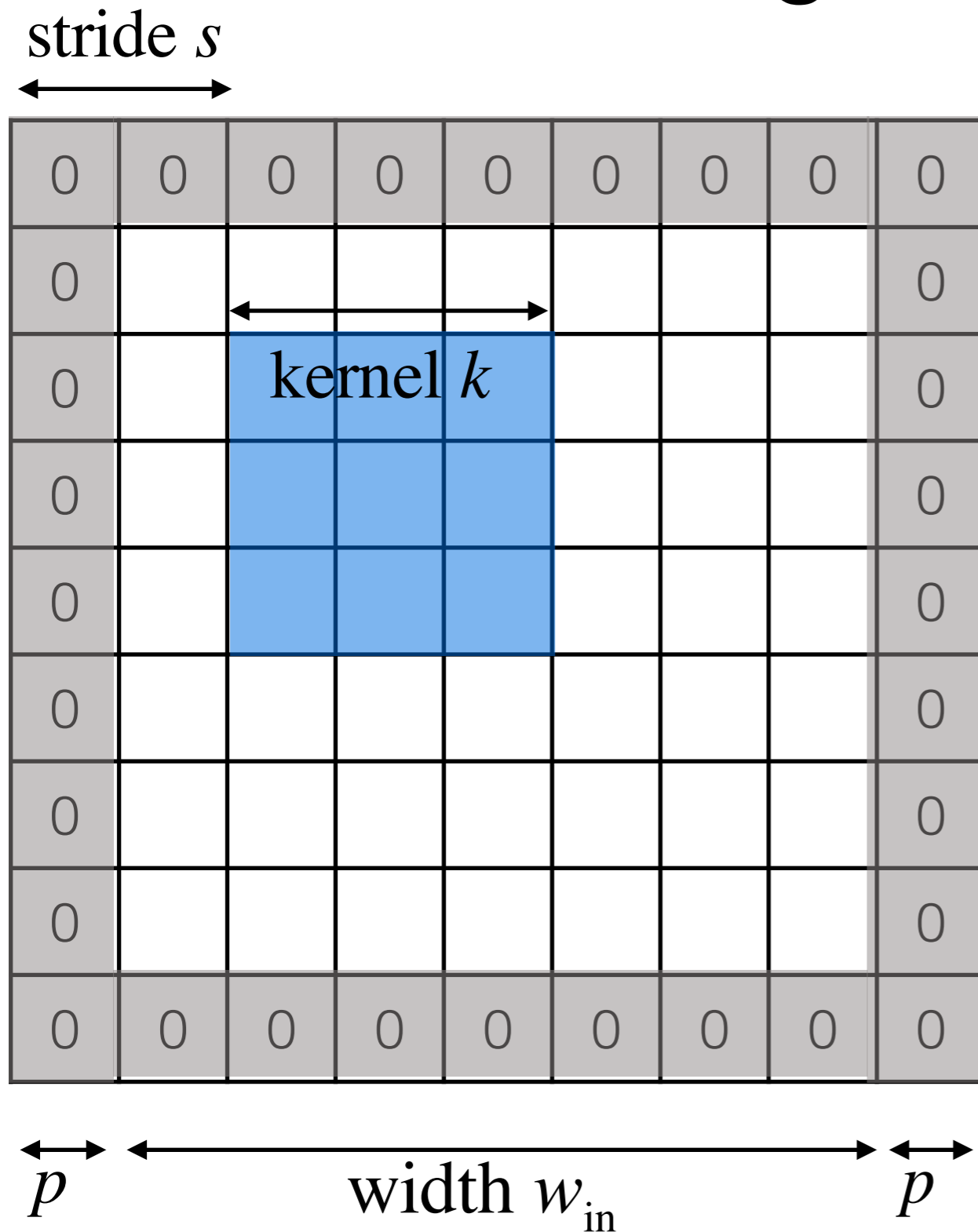
How big is the output?



Example: $k=3$, $s=1$, $p=1$

Convolution:

How big is the output?

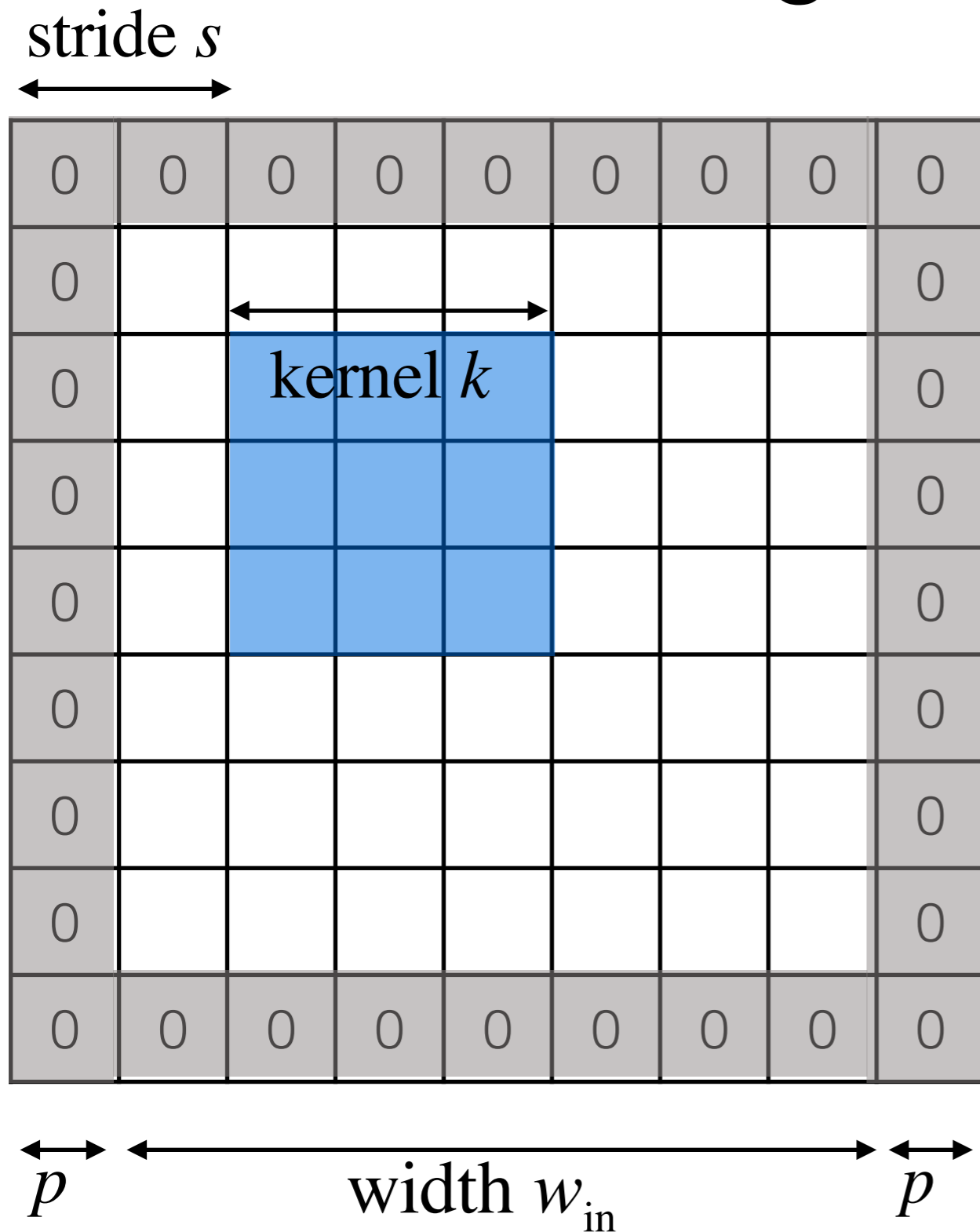


Example: $k=3$, $s=1$, $p=1$

$$\begin{aligned}w_{\text{out}} &= \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1 \\ &= \left\lfloor \frac{w_{\text{in}} + 2 - 3}{1} \right\rfloor + 1 \\ &= w_{\text{in}}\end{aligned}$$

Convolution:

How big is the output?



Example: $k=3$, $s=1$, $p=1$

$$\begin{aligned}w_{\text{out}} &= \left\lfloor \frac{w_{\text{in}} + 2p - k}{s} \right\rfloor + 1 \\ &= \left\lfloor \frac{w_{\text{in}} + 2 - 3}{1} \right\rfloor + 1 \\ &= w_{\text{in}}\end{aligned}$$

VGGNet [Simonyan 2014]
uses filters of this shape

Max Pooling

For most CNNs, **convolution** is often followed by **pooling**:

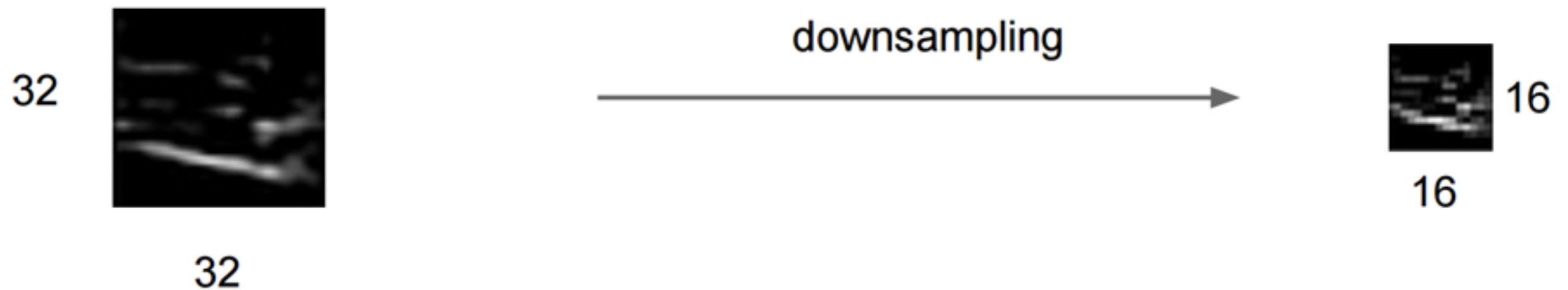
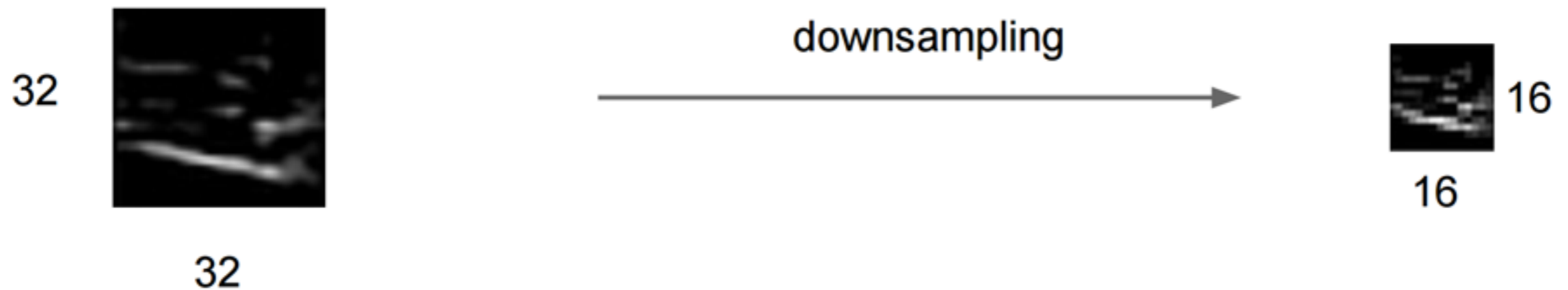


Figure: Andrej Karpathy

Max Pooling

For most CNNs, **convolution** is often followed by **pooling**:

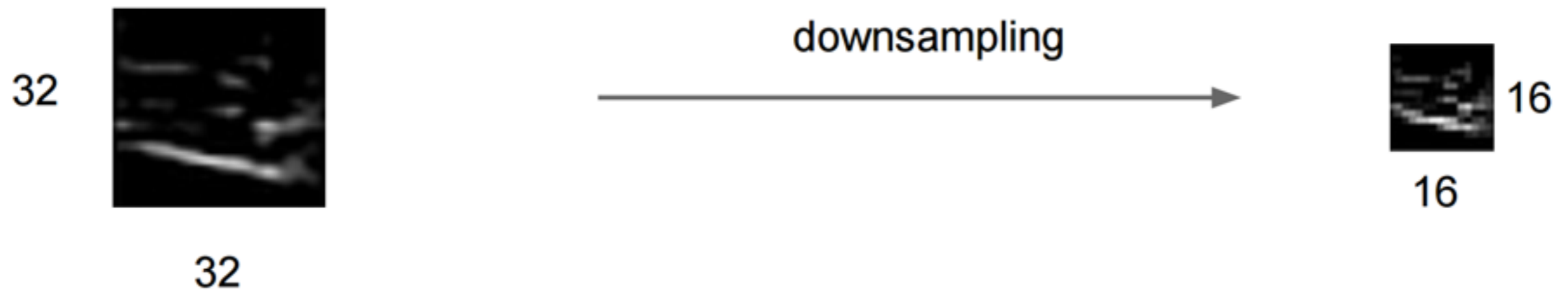
- Creates a smaller representation while retaining the most important information



Max Pooling

For most CNNs, **convolution** is often followed by **pooling**:

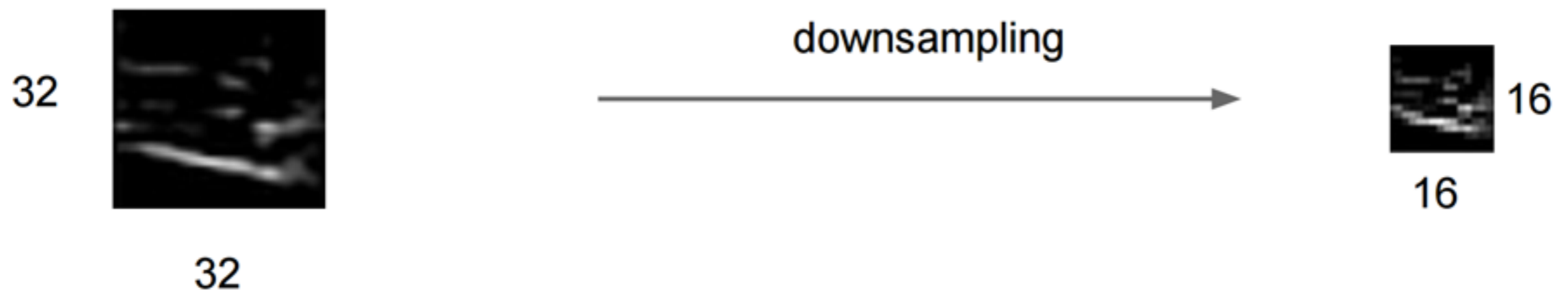
- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common



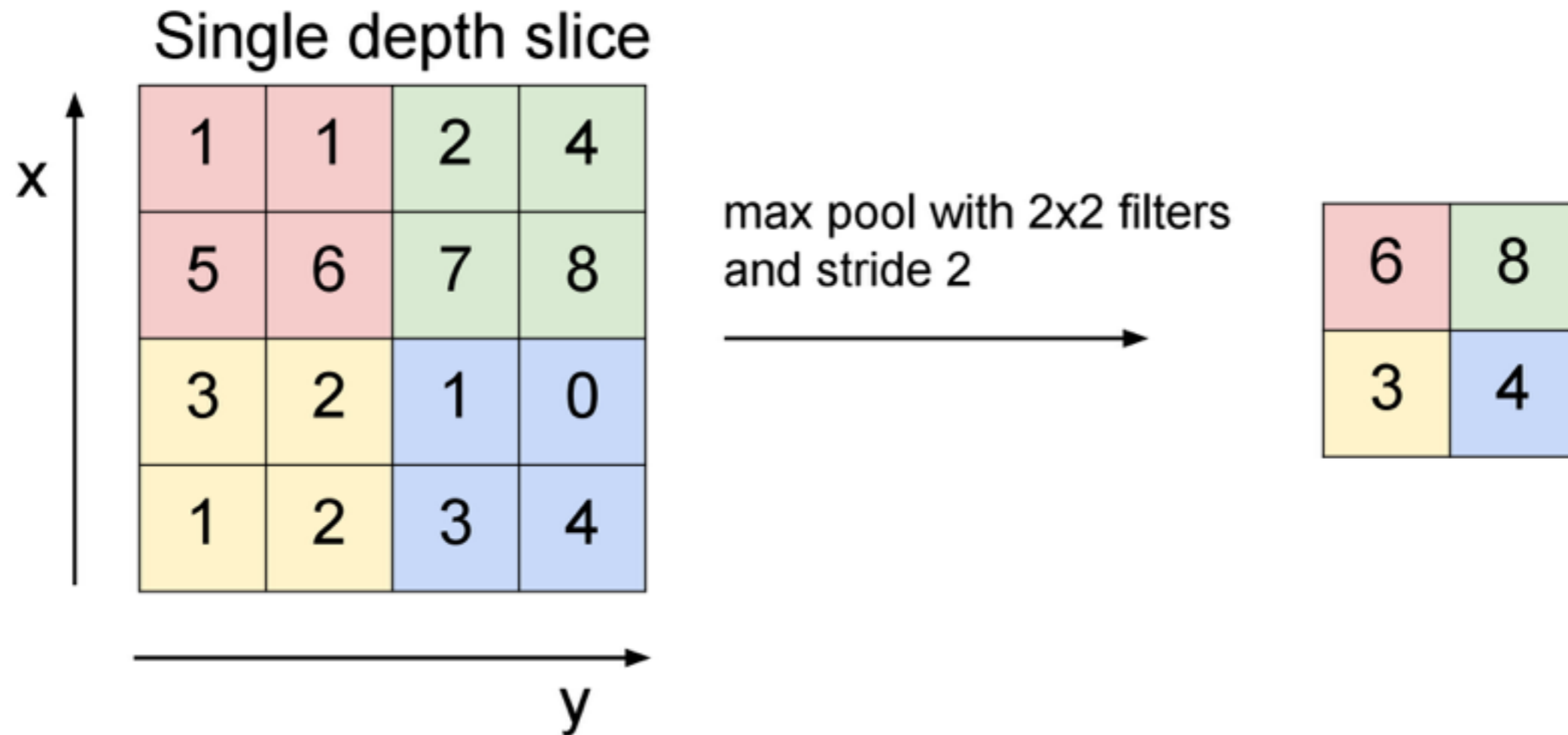
Max Pooling

For most CNNs, **convolution** is often followed by **pooling**:

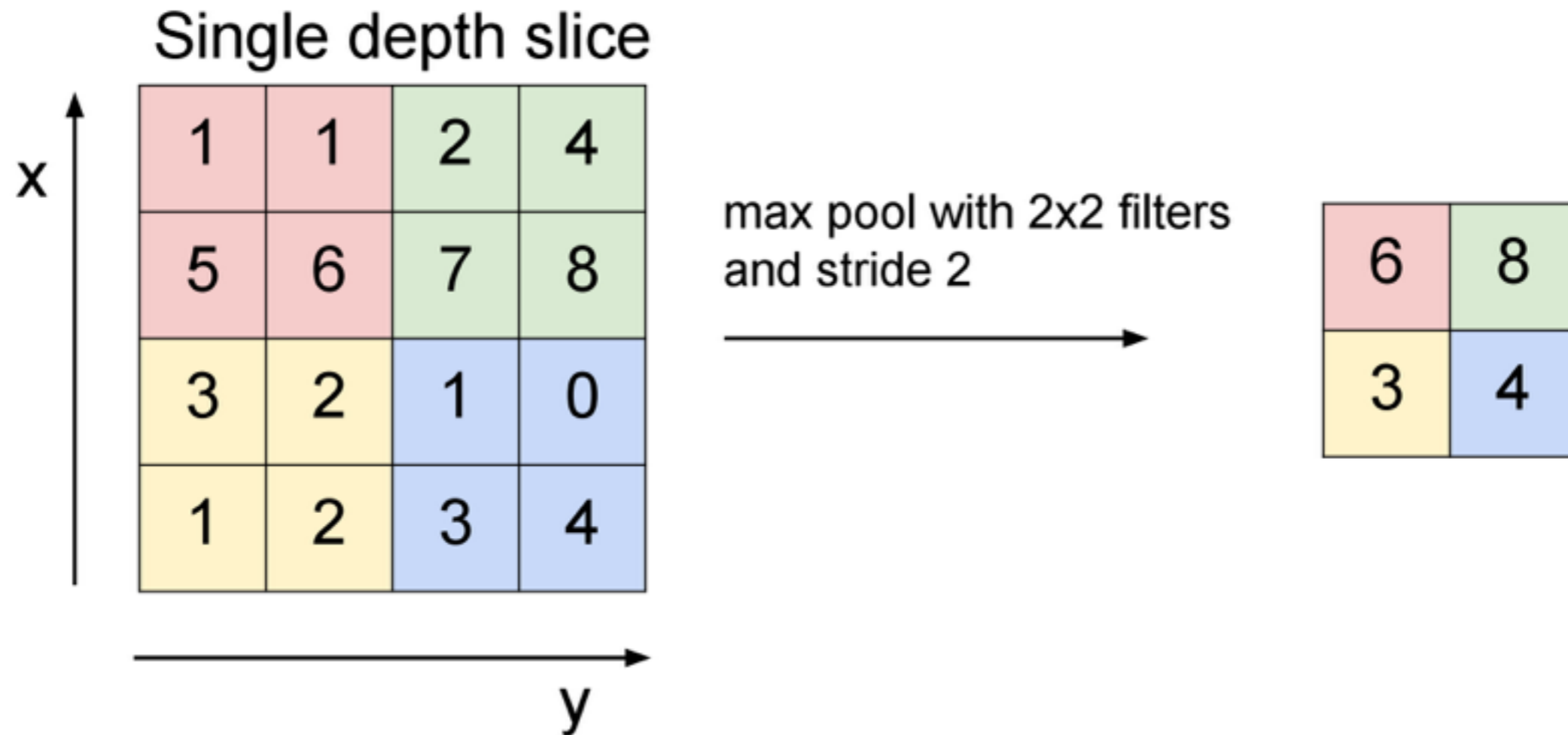
- Creates a smaller representation while retaining the most important information
- The “max” operation is the most common
- Why might “avg” be a poor choice?



Max Pooling

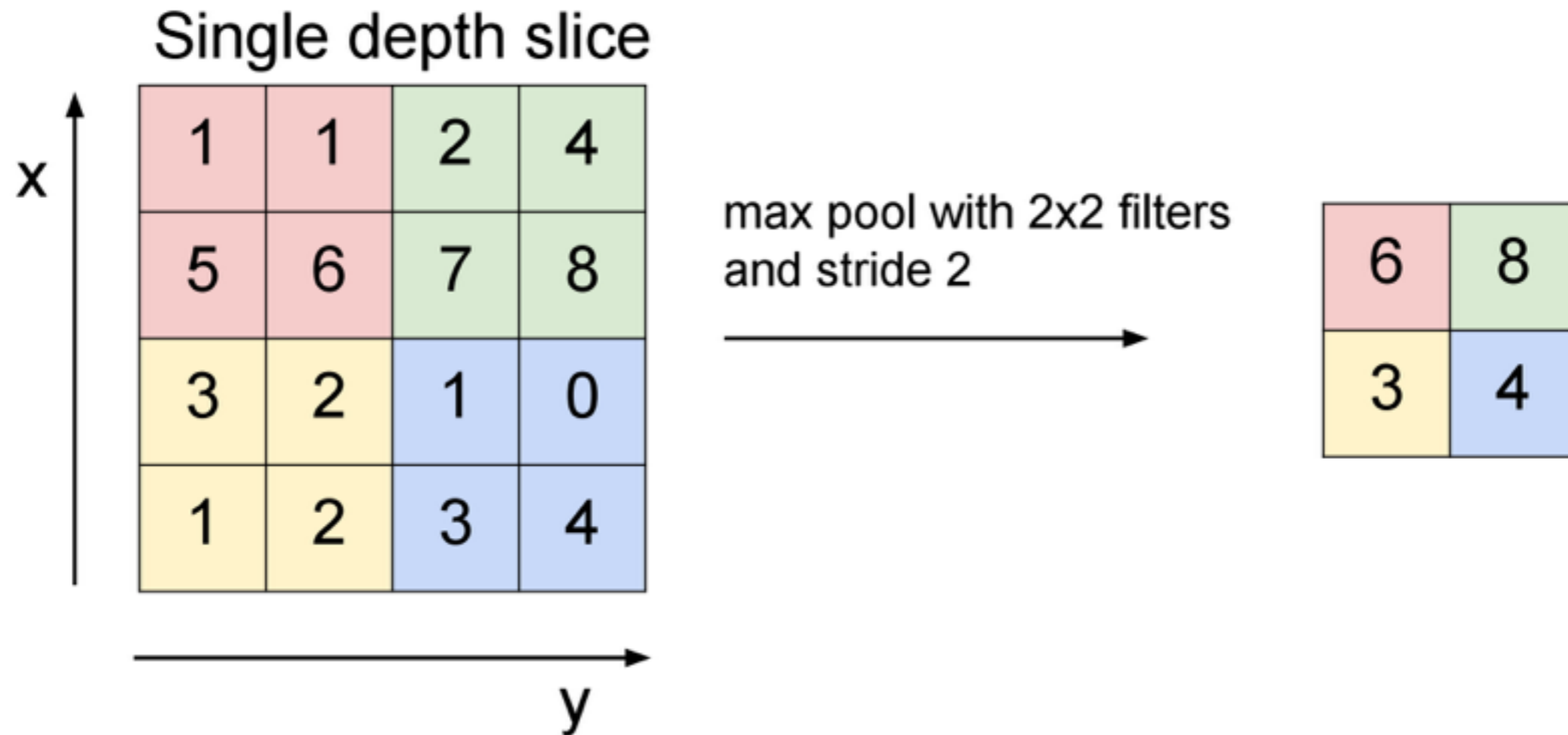


Max Pooling



What's the backprop rule for max pooling?

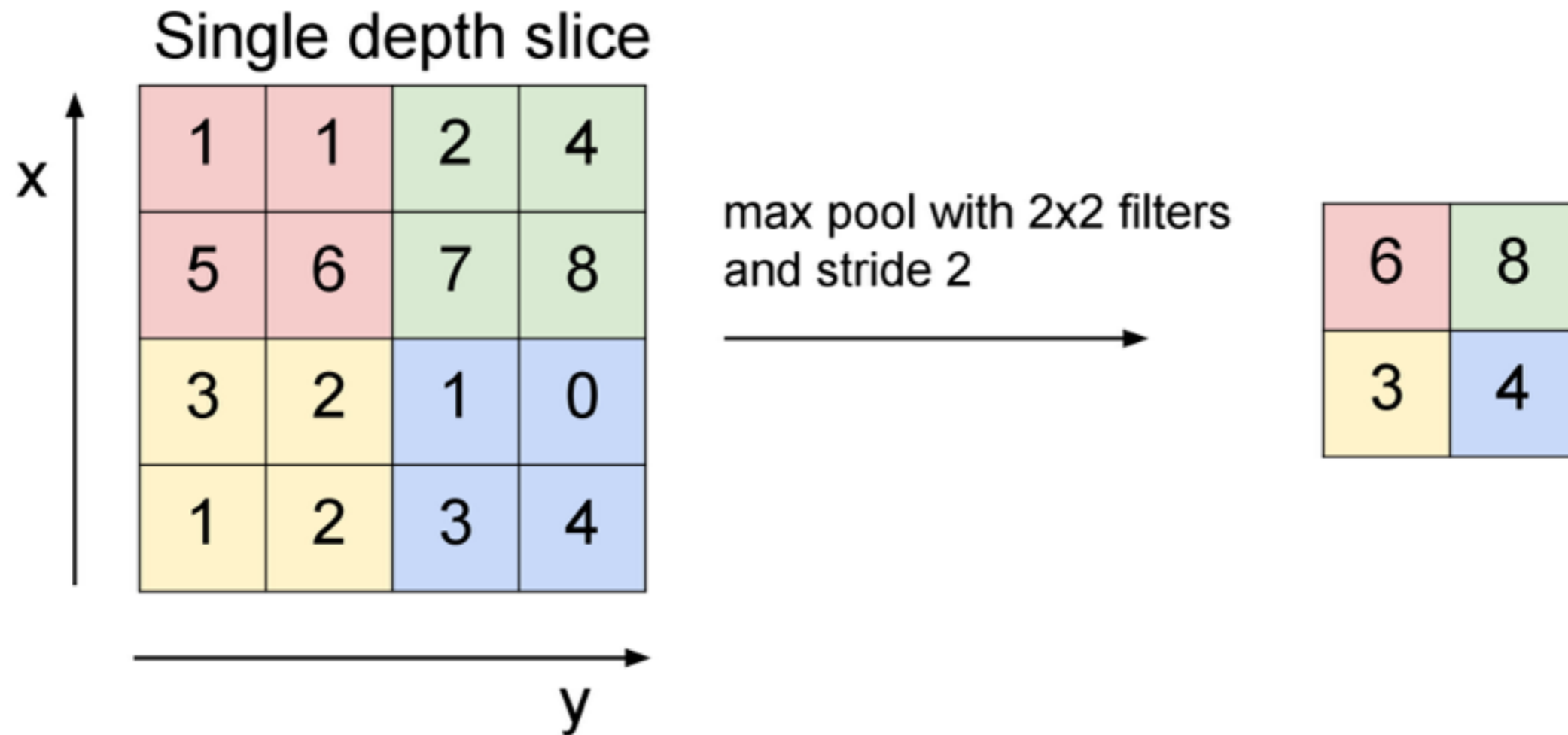
Max Pooling



What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max

Max Pooling



What's the backprop rule for max pooling?

- In the forward pass, store the index that took the max
- The backprop gradient is the input gradient at that index

Example CNN



Figure: Andrej Karpathy

Example CNN

CONV CONV POOL CONV CONV POOL CONV CONV POOL
↓ ReLU ↓ ReLU ↓ ↓ ReLU ↓ ReLU ↓ ↓ ReLU ↓ ReLU ↓

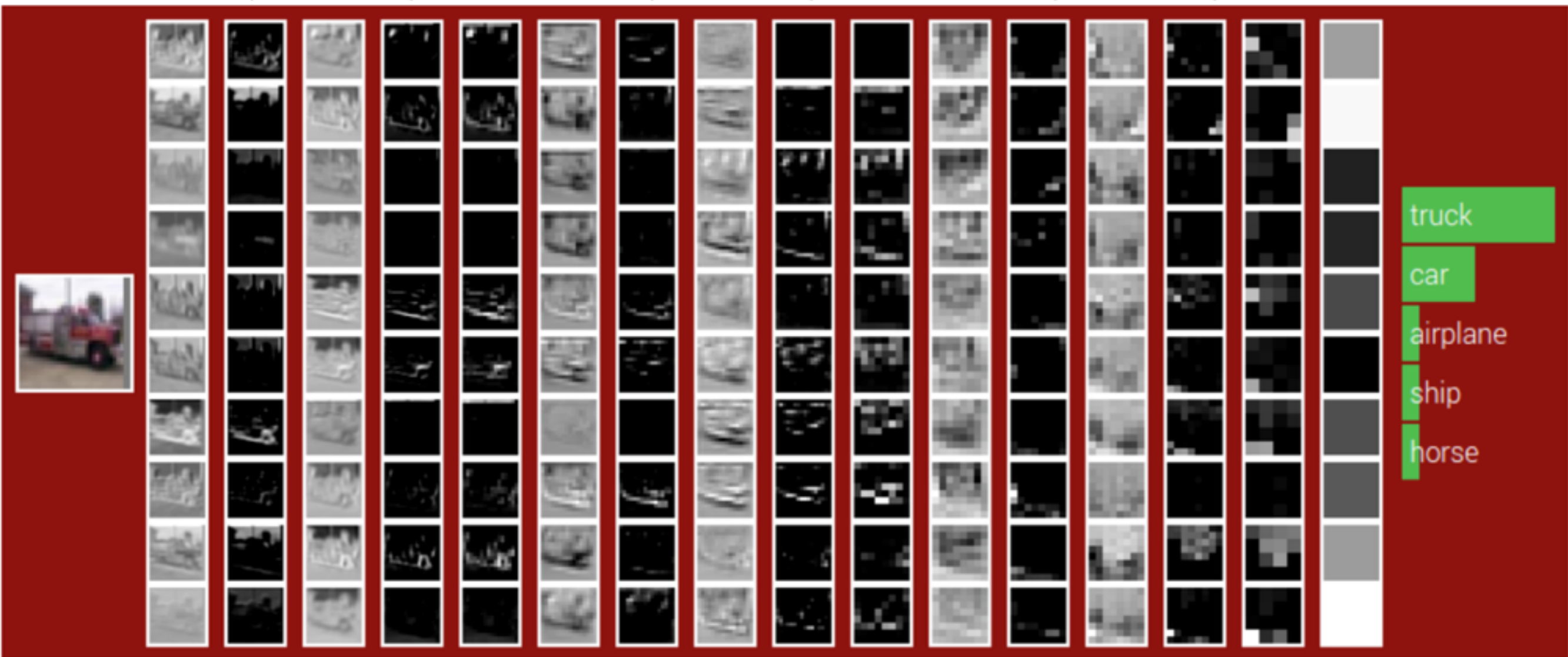


Figure: Andrej Karpathy

Example CNN

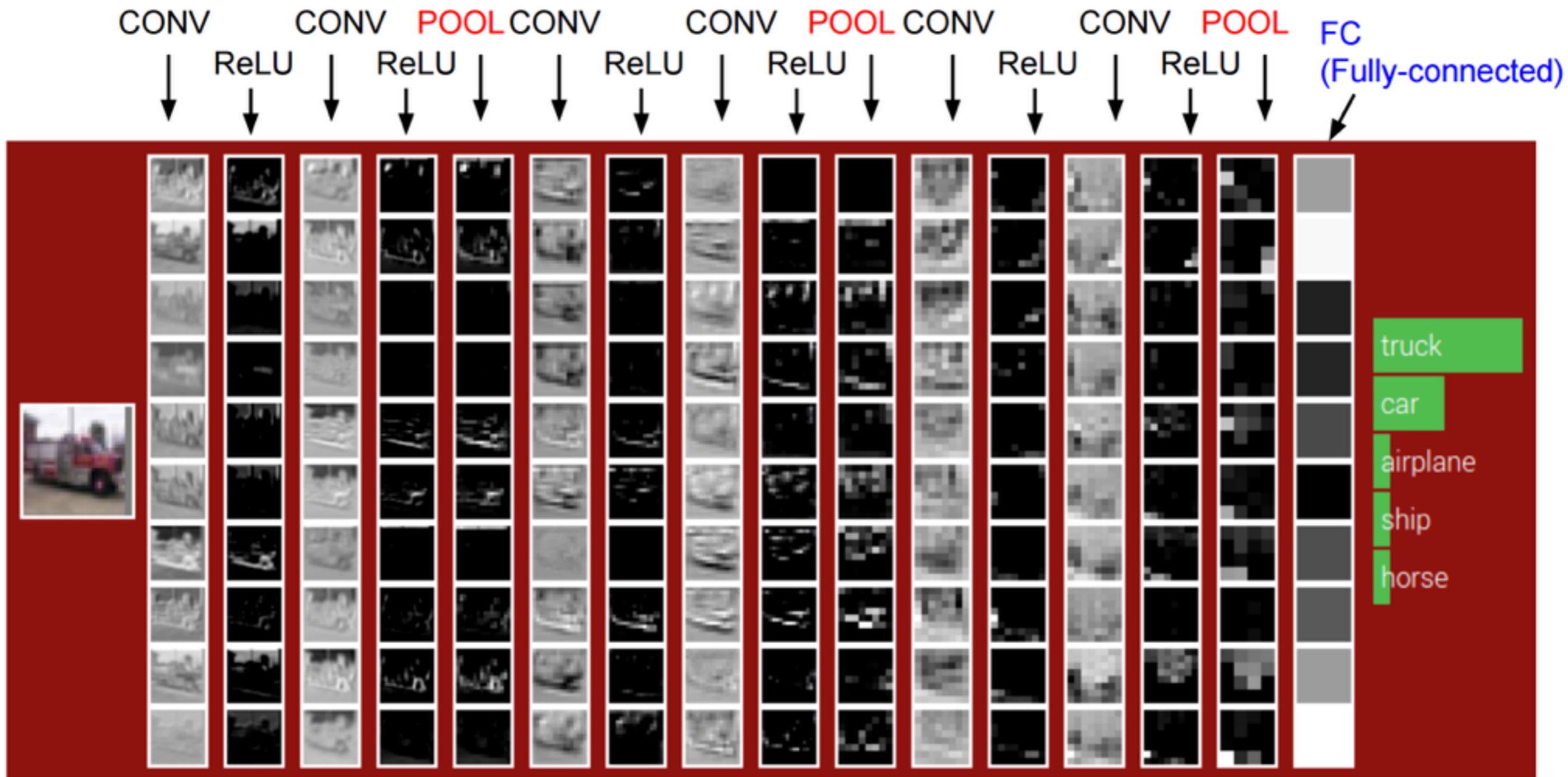
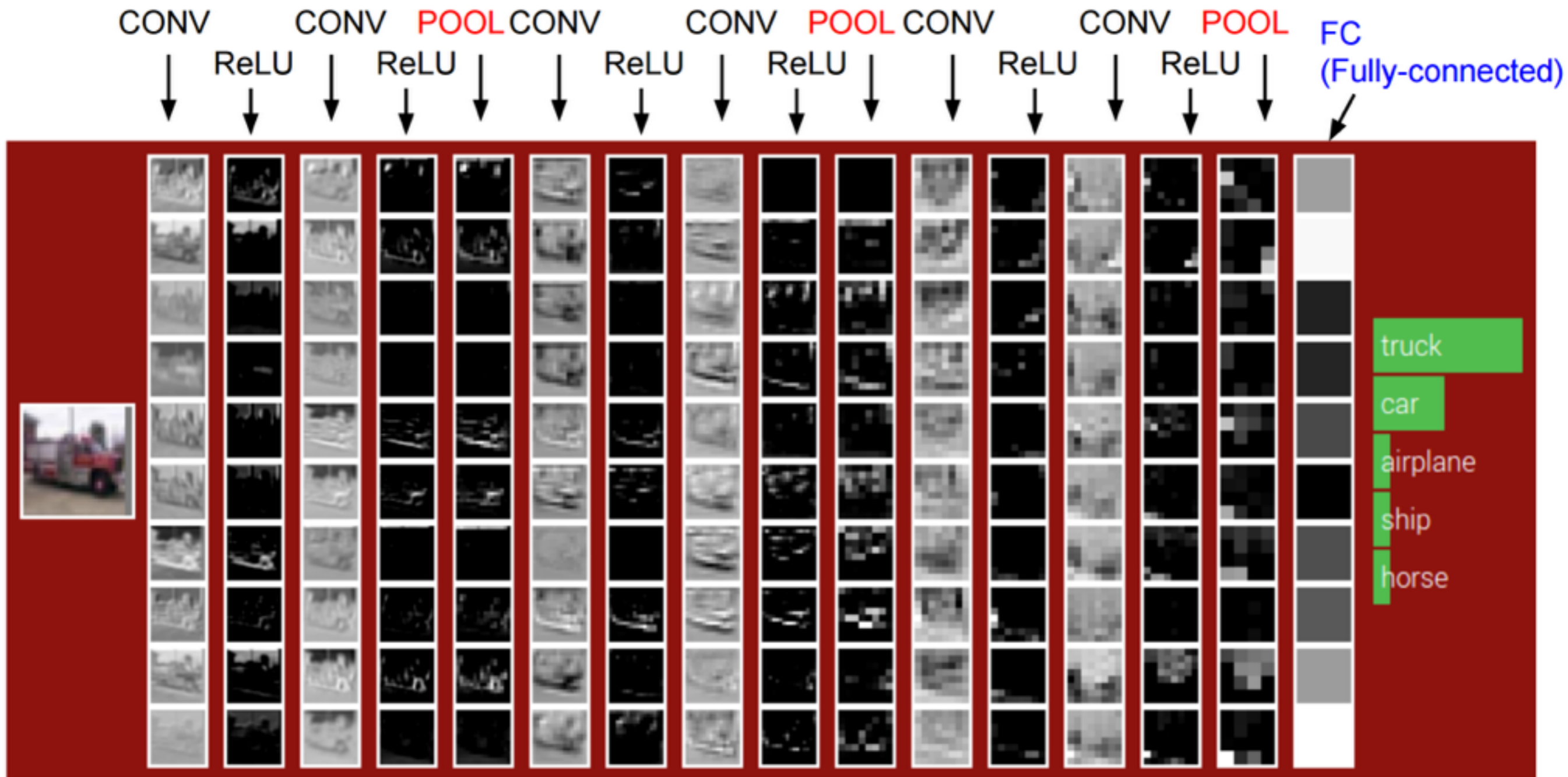


Figure: Andrej Karpathy

Example CNN



10x3x3 conv filters, stride 1, pad 1
2x2 pool filters, stride 2

Figure: Andrej Karpathy

Questions?