

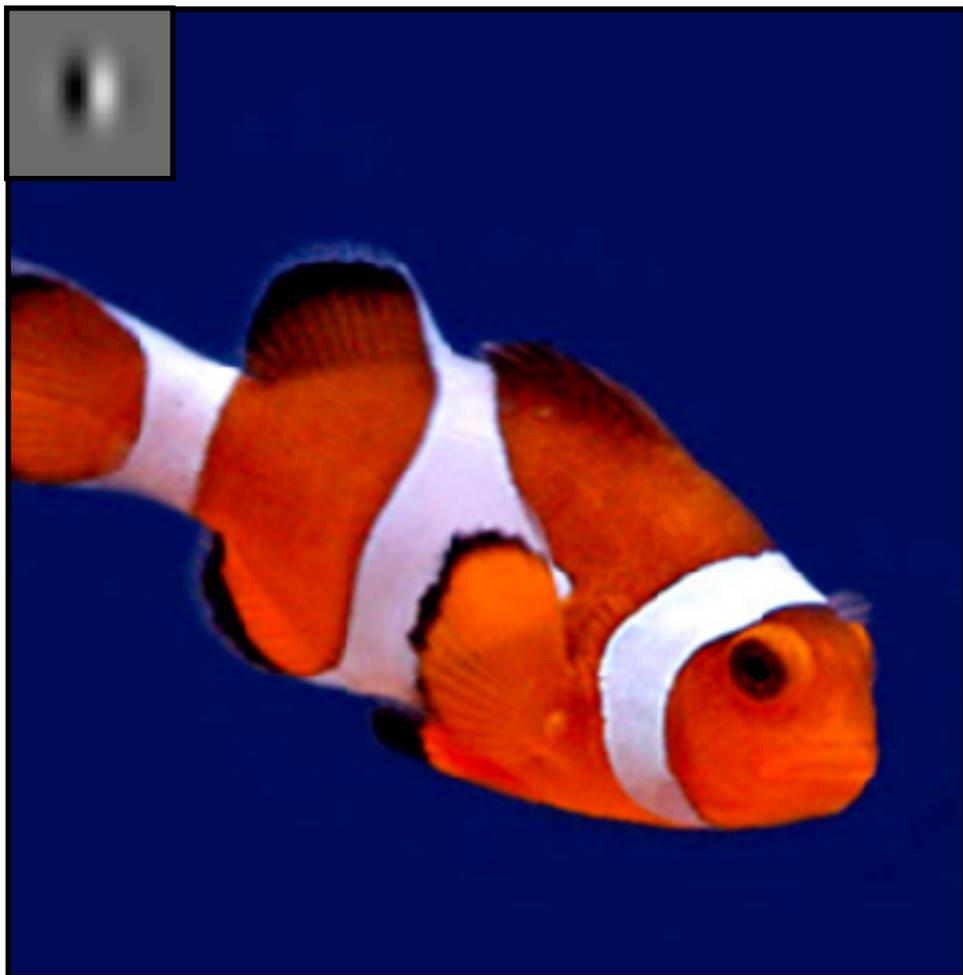
Deep Learning for Analyzing Images and Time Series

nearly all slides by George Chen (CMU)

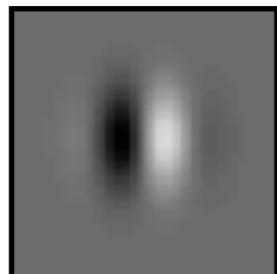
1 slide by Phillip Isola (OpenAI, UC Berkeley)

Image Analysis with Convolutional Neural Nets (CNNs, also called convnets)

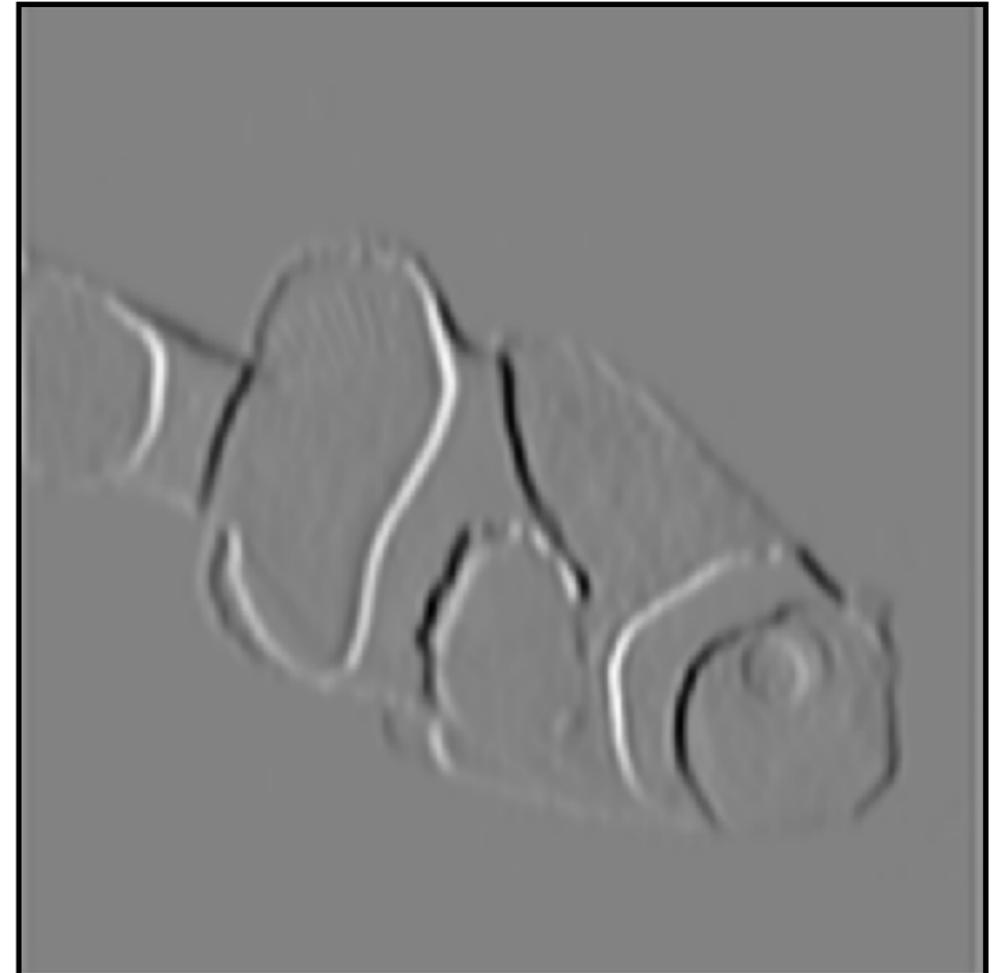
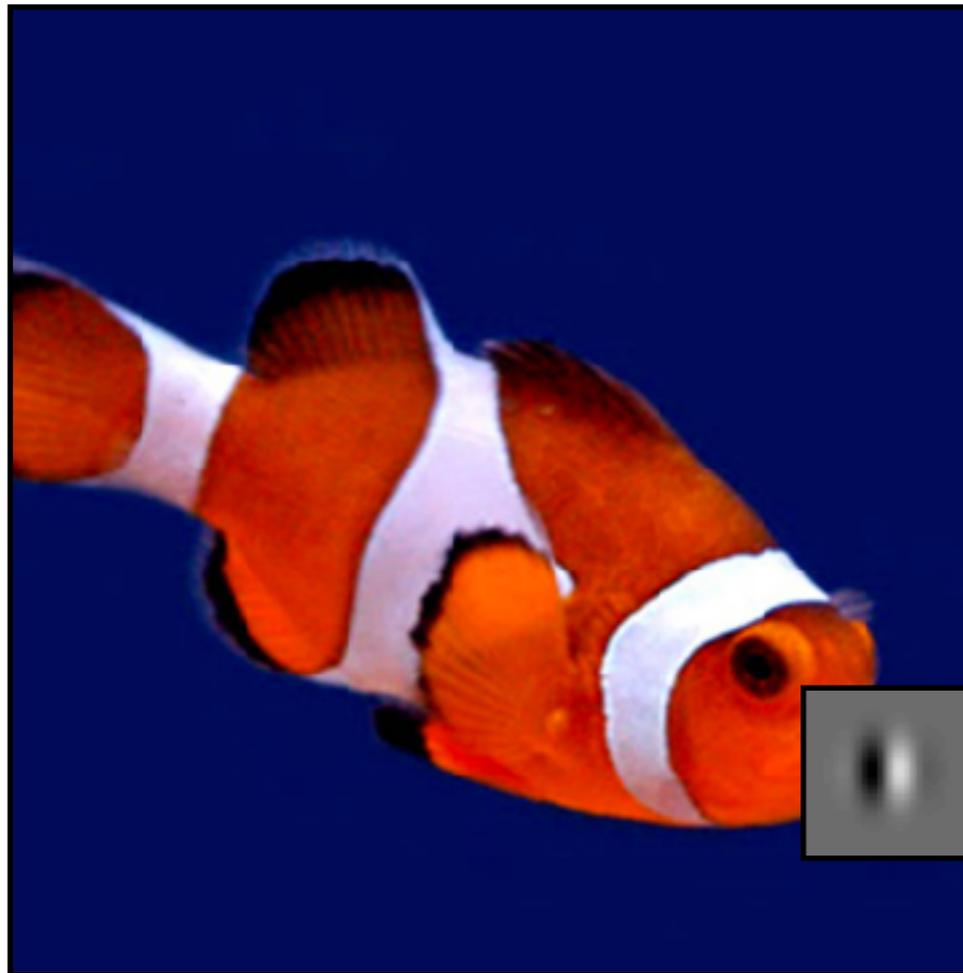
Convolution



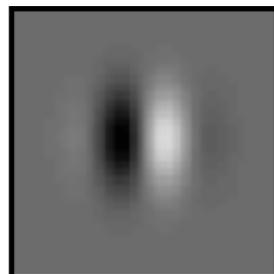
filter



Convolution



filter



Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	0	0
0	1	0
0	0	0

Filter
(also called "kernel")

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	0	0
0	1	0
0	0	0

Filter
(also called "kernel")

Convolution

0	0	0	0	0	0	0
0	0	1	0	1	1	0
0	1	0	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	0	1	1	0
0	1	0	0	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	0	1	1	0
0	1	0	0	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	0	1	1	0
0	1	0	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0				

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	1			

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	0	1	0	1
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	1	1		

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	1	1	1	

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	0	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

0	1	1	1	0

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0		
0	0	0	1	0	1	1	0	0
0	0	1	1	1	0	1	1	0
0	0	1	0	1	0	1	0	0
0	1	1	1	1	1	1	0	
0	0	1	1	1	0	0		
0	0	0	0	0	0	0		

Input image

0	1	1	1	0
1				

Output image

Convolution

Take dot product!

0	0	0	0	0	0	0	
0	0	1	1	0	1	0	0
0	1	0	1	1	0	1	0
0	1	0	1	0	0	0	0
0	1	1	1	1	1	1	0
0	0	1	1	1	0	0	0
0	0	0	0	0	0	0	0

Input image

0	1	1	1	0
1	1			

Output image

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

*

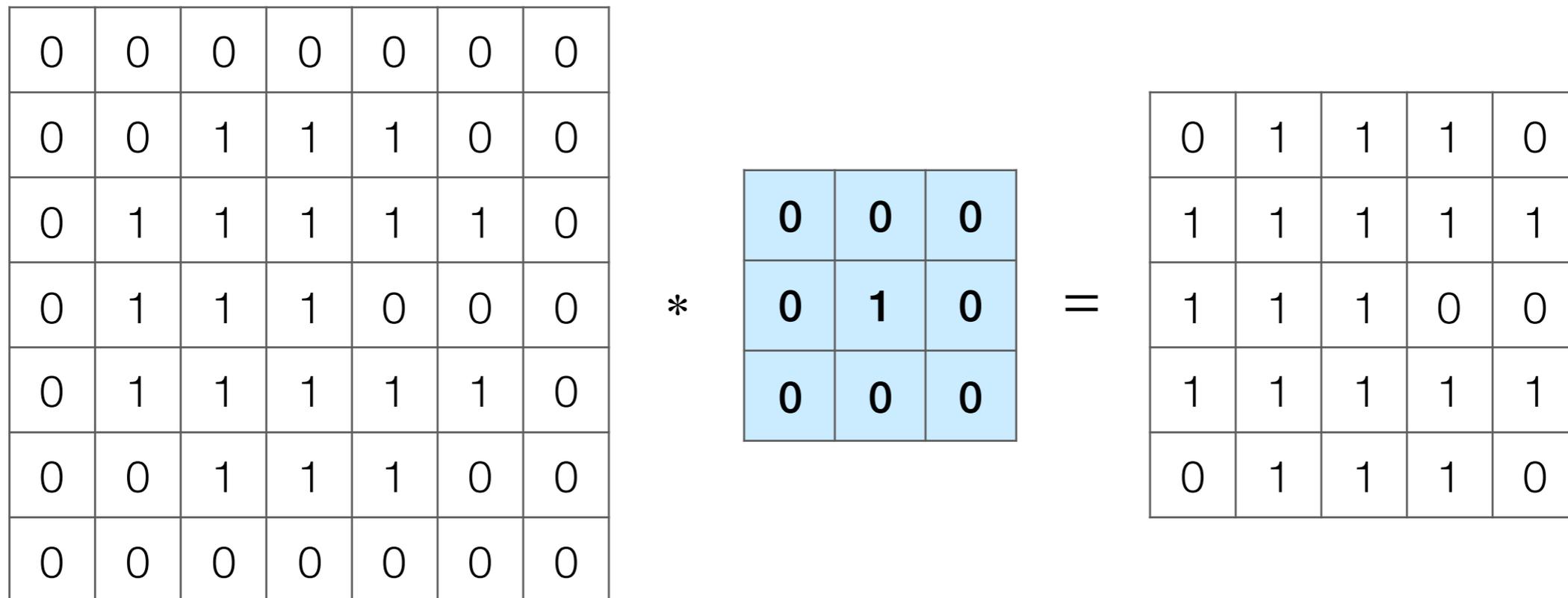
0	0	0
0	1	0
0	0	0

=

0	1	1	1	0
1	1	1	1	1
1	1	1	0	0
1	1	1	1	1
0	1	1	1	0

Output image

Convolution



Input image

Output image

Note: output image is smaller than input image

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

*

0	0	0
0	1	0
0	0	0

=

0	1	1	1	0
1	1	1	1	1
1	1	1	0	0
1	1	1	1	1
0	1	1	1	0

Input image

Output image

Note: output image is smaller than input image

If you want output size to be same as input, pad 0's to input

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

*

0	0	0
0	1	0
0	0	0

=

0	1	1	1	0
1	1	1	1	1
1	1	1	0	0
1	1	1	1	1
0	1	1	1	0

Input image

Output image

Note: output image is smaller than input image

If you want output size to be same as input, pad 0's to input

Convolution

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	1	1	1	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	0	1	1	1	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Input image

*

0	0	0
0	1	0
0	0	0

=

0	1	1	1	0
1	1	1	1	1
1	1	1	0	0
1	1	1	1	1
0	1	1	1	0

Output image

Note: output image is smaller than input image

If you want output size to be same as input, pad 0's to input

Convolution

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	1	1	1	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	0	1	1	1	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Input image

*

0	0	0
0	1	0
0	0	0

=

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Output image

Note: output image is smaller than input image

If you want output size to be same as input, pad 0's to input

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

*

0	0	0
0	1	0
0	0	0

=

0	1	1	1	0
1	1	1	1	1
1	1	1	0	0
1	1	1	1	1
0	1	1	1	0

Output image

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

*

1	1	1
1	1	1
1	1	1

=

3	5	6	5	3
5	8	8	6	3
6	9	8	7	4
5	8	8	6	3
3	5	6	5	3

Output image

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

*	$\frac{1}{9}$	1	1	1
		1	1	1
		1	1	1

=	$\frac{1}{9}$	3	5	6	5	3
		5	8	8	6	3
		6	9	8	7	4
		5	8	8	6	3
		3	5	6	5	3

Output image

Convolution

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

*

-1	-1	-1
2	2	2
-1	-1	-1

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

Output image

Convolution

Convolution

Very commonly used for:

Convolution

Very commonly used for:

- Blurring an image

Convolution

Very commonly used for:

- Blurring an image



*

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

=



Convolution

Very commonly used for:

- Blurring an image



$$\begin{matrix} * & \begin{matrix} \begin{matrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{matrix} \\ = \end{matrix} \end{matrix}$$



- Finding edges

Convolution

Very commonly used for:

- Blurring an image



$$\begin{matrix} * & \begin{matrix} \begin{matrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{matrix} \\ & = \end{matrix} \end{matrix}$$



- Finding edges



$$\begin{matrix} * & \begin{matrix} \begin{matrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{matrix} \\ & = \end{matrix} \end{matrix}$$



(this example finds horizontal edges)

Convolution Layer

Convolution Layer



Convolution Layer



1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

-1	-1	-1
2	2	2
-1	-1	-1

0	-1	0
-1	4	-1
0	-1	0

Convolution Layer



convolve with
each filter

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

-1	-1	-1
2	2	2
-1	-1	-1

0	-1	0
-1	4	-1
0	-1	0

Convolution Layer



1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

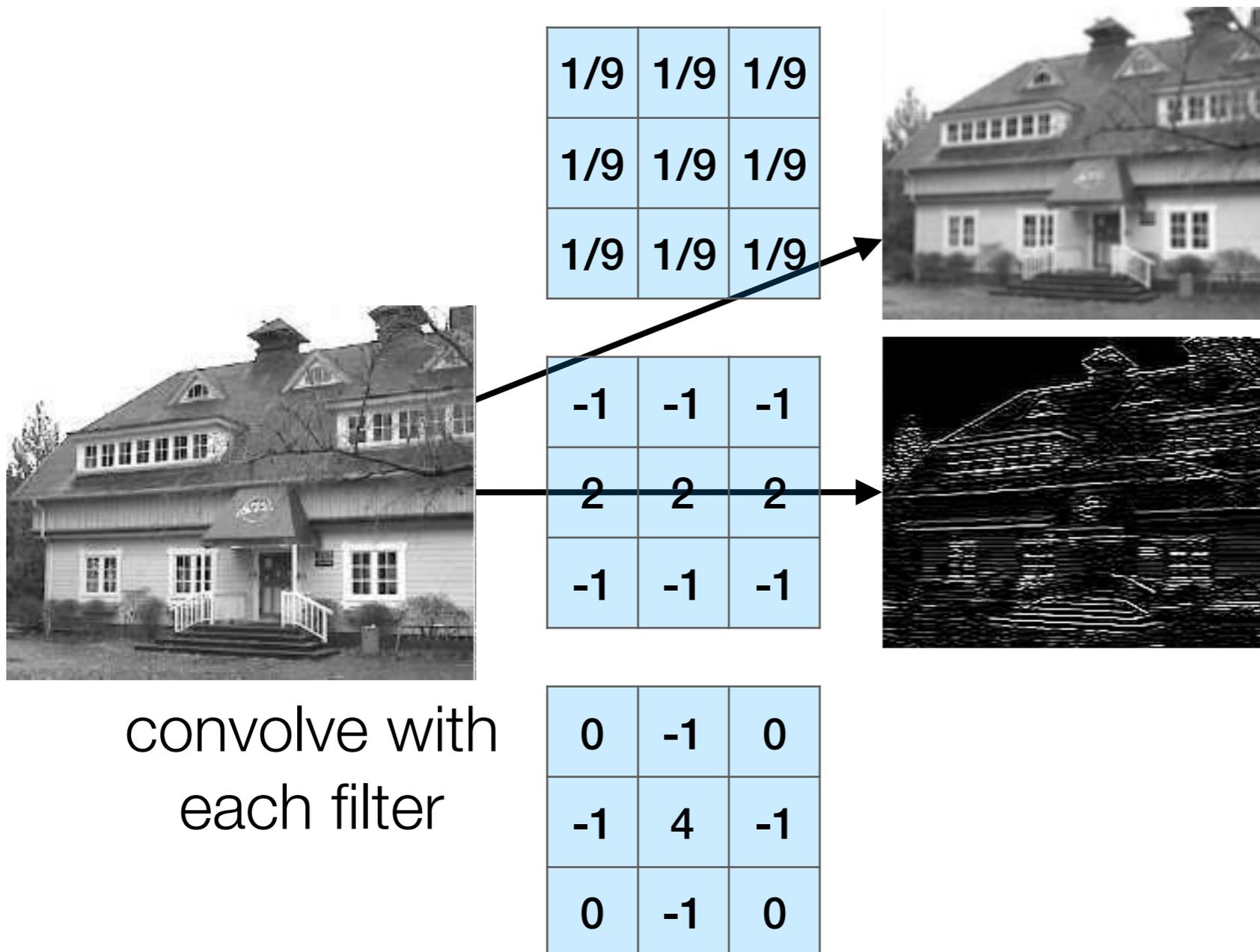


-1	-1	-1
2	2	2
-1	-1	-1

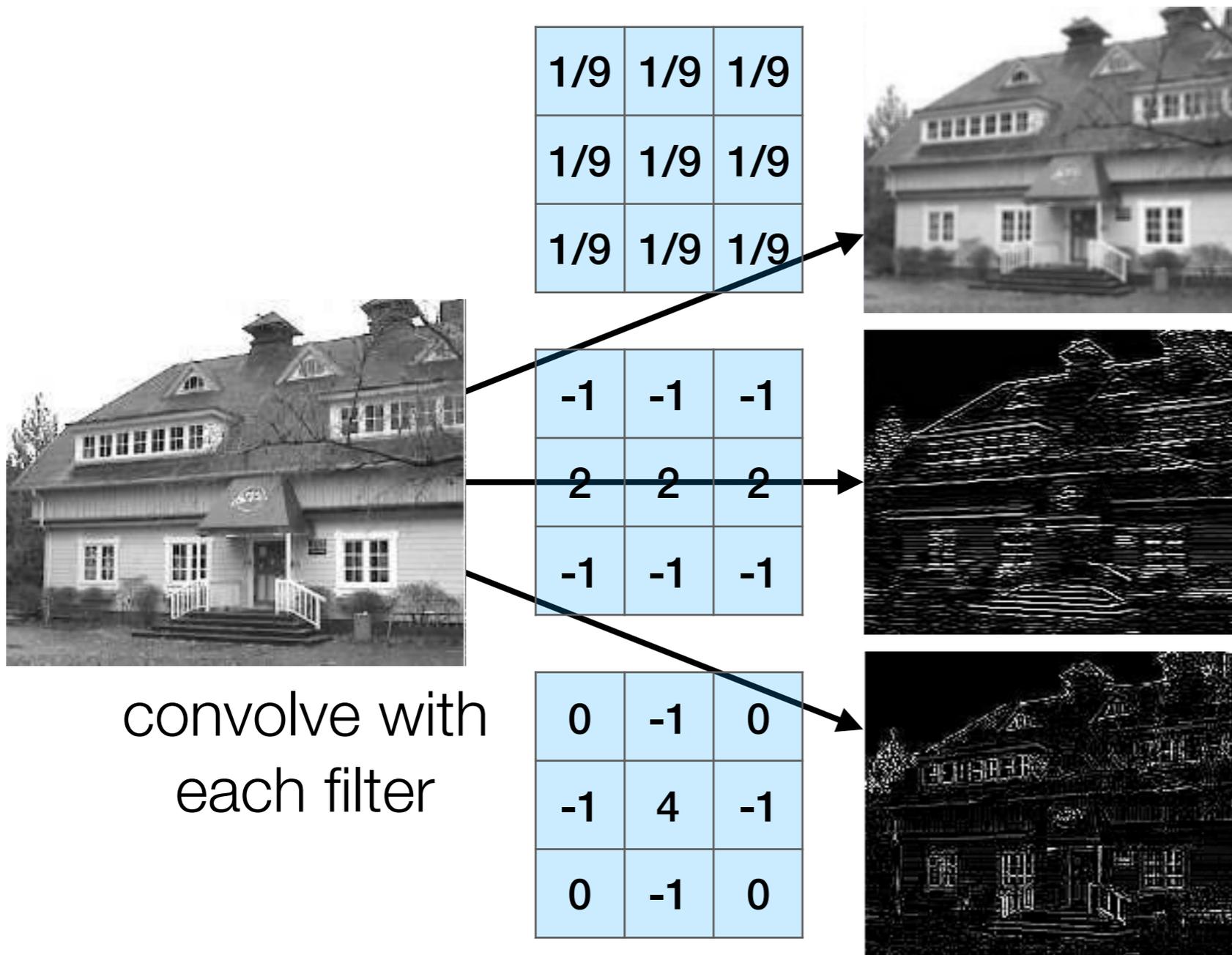
convolve with
each filter

0	-1	0
-1	4	-1
0	-1	0

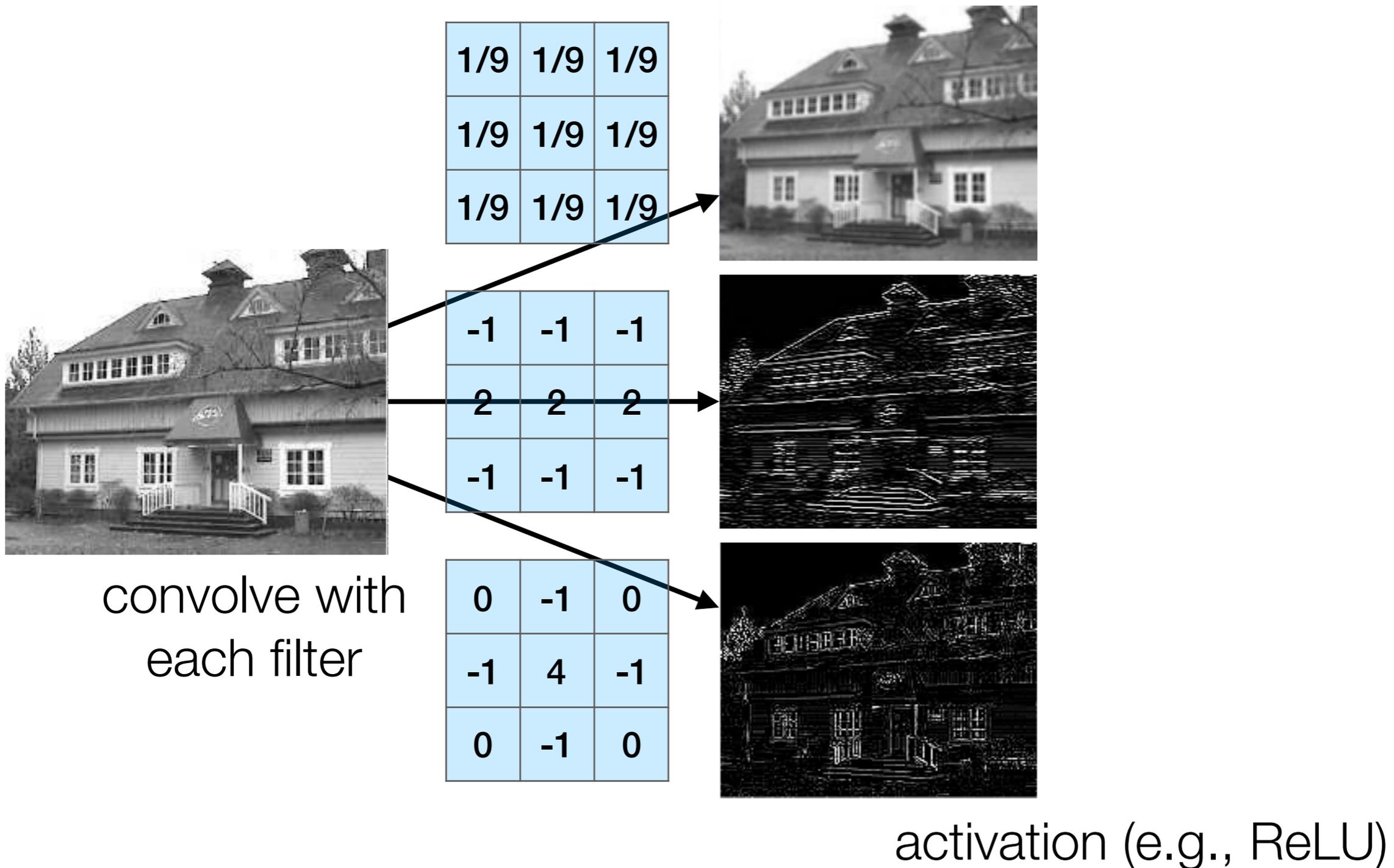
Convolution Layer



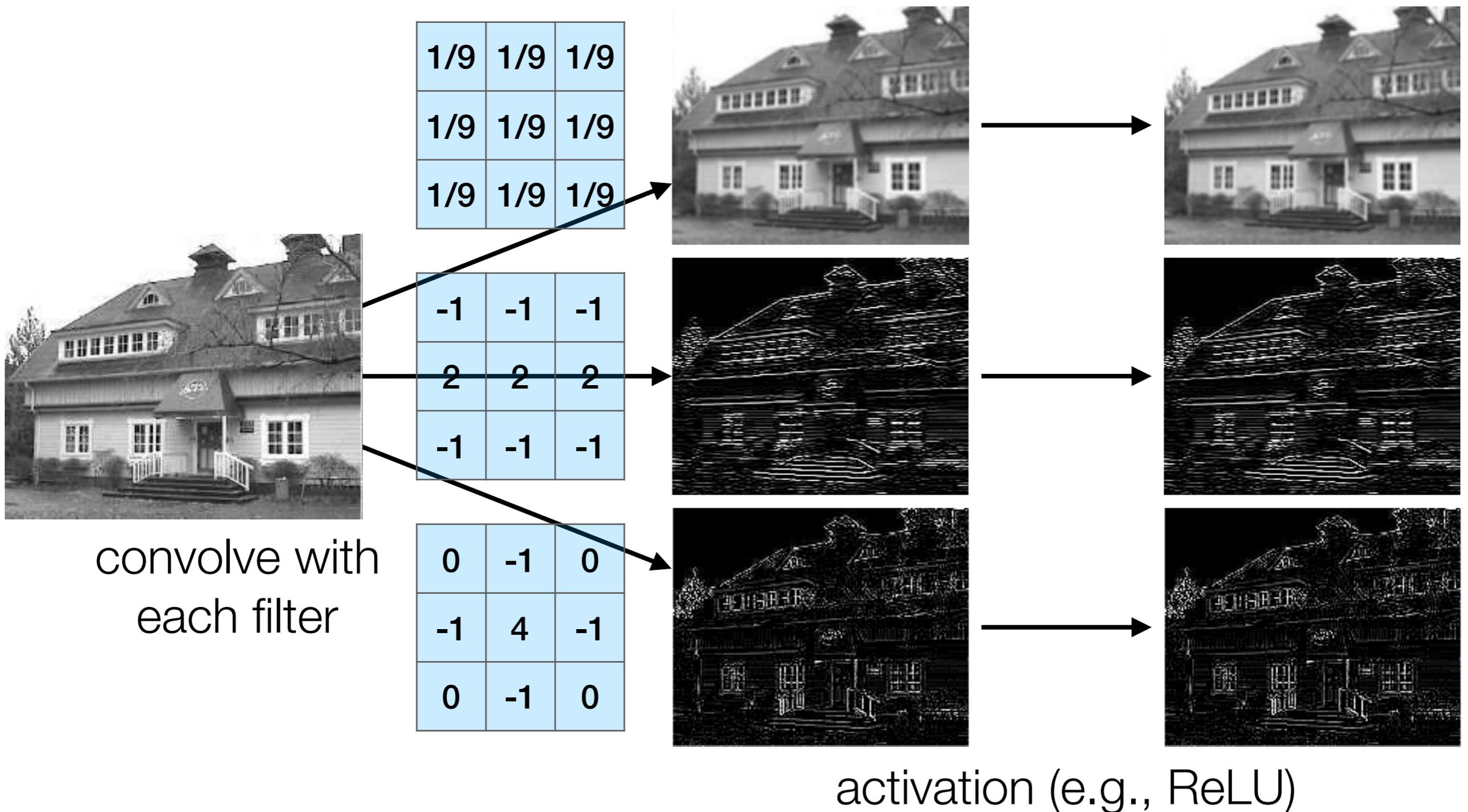
Convolution Layer



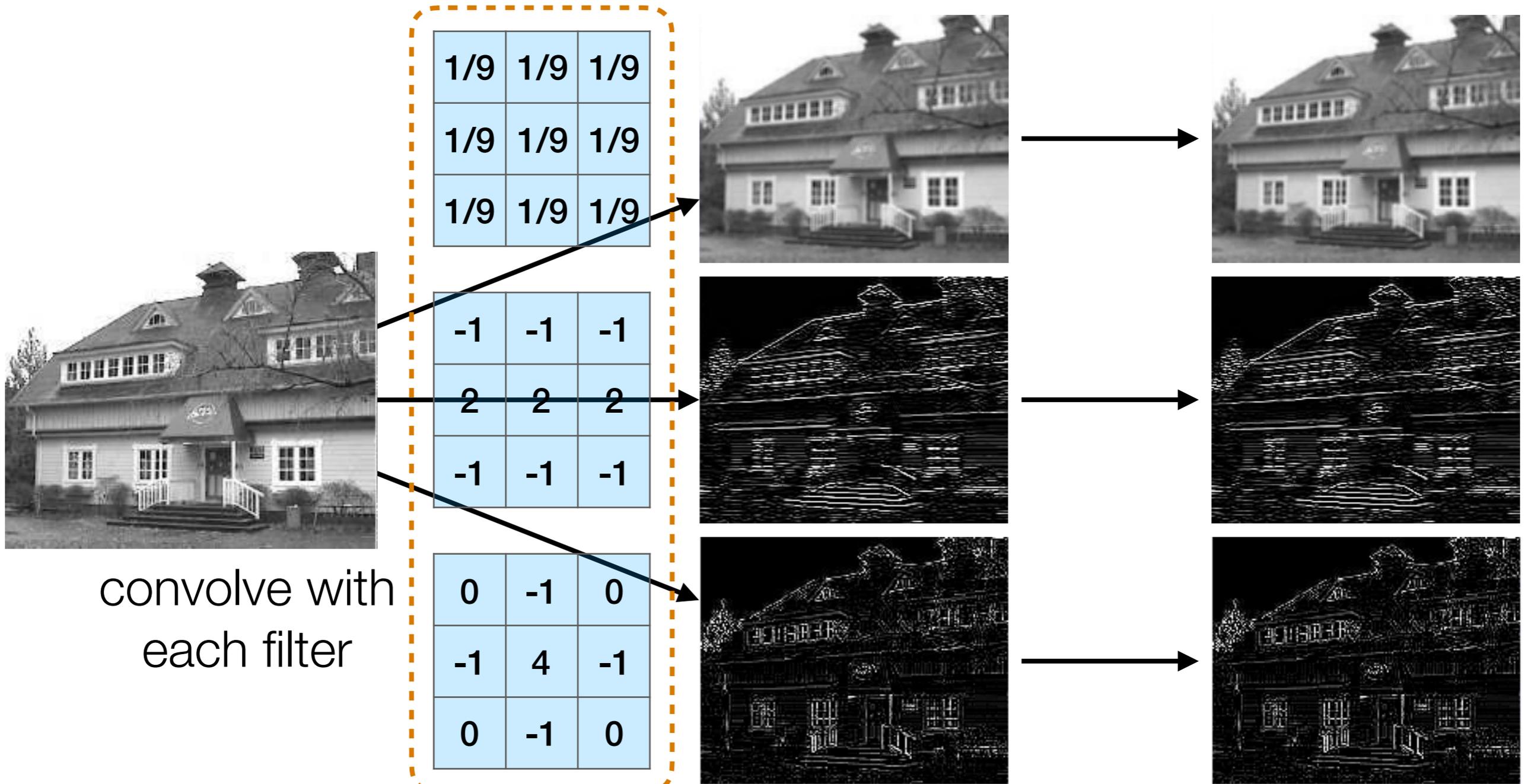
Convolution Layer



Convolution Layer



Convolution Layer

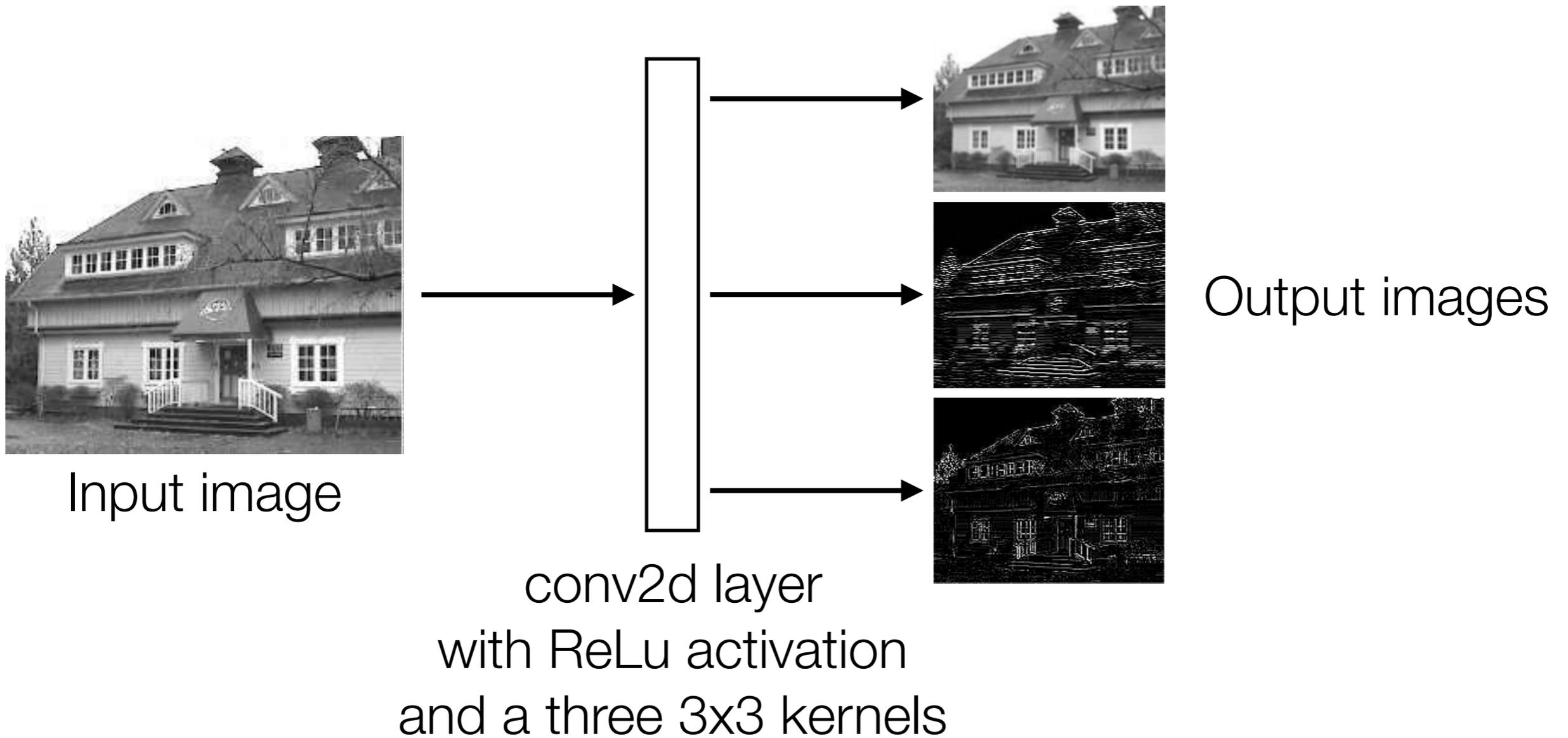


convolve with each filter

filters are actually unknown and are learned!

activation (e.g., ReLU)

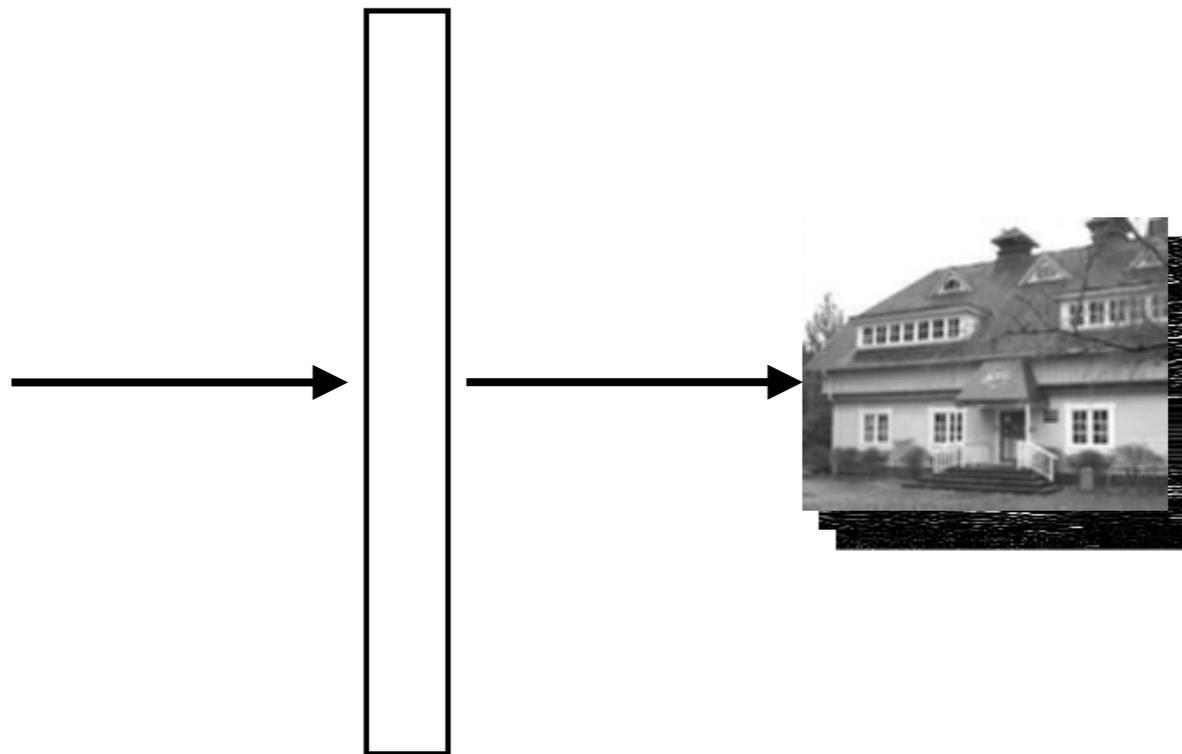
Convolution Layer



Convolution Layer



Input image



conv2d layer
with ReLu activation
and a three 3x3 kernels



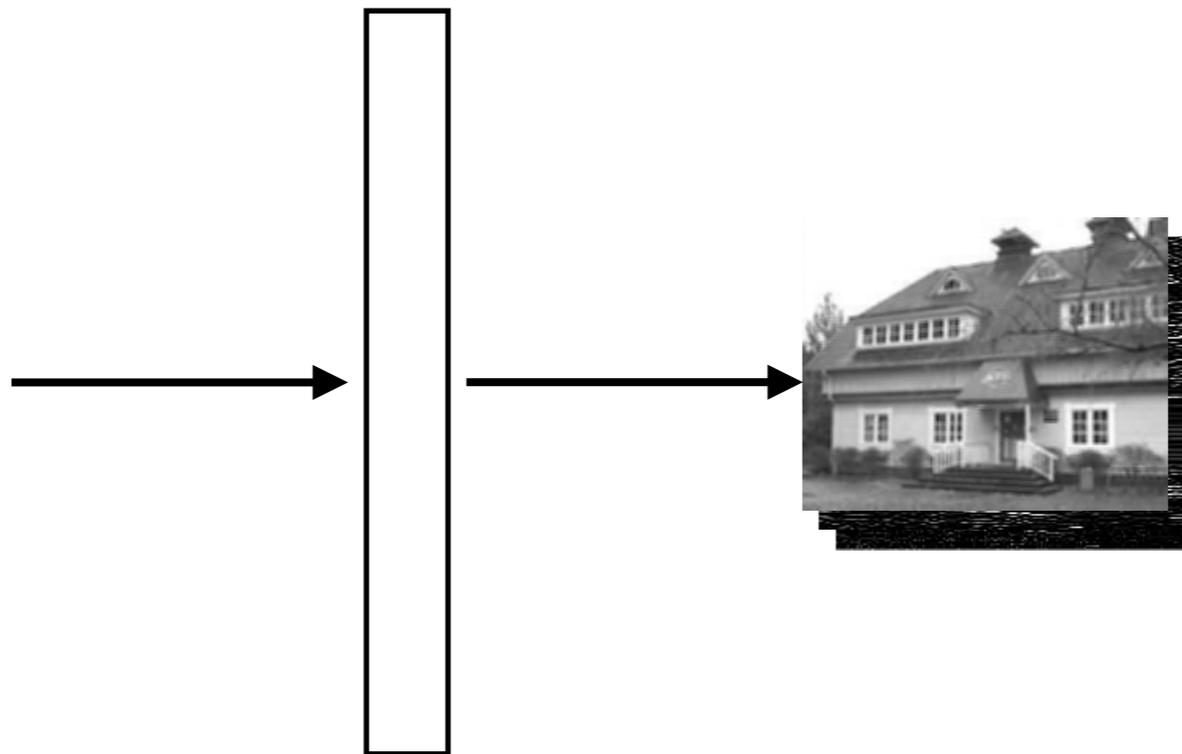
Stack output
images into a
single “output
feature map”

Convolution Layer



Input image

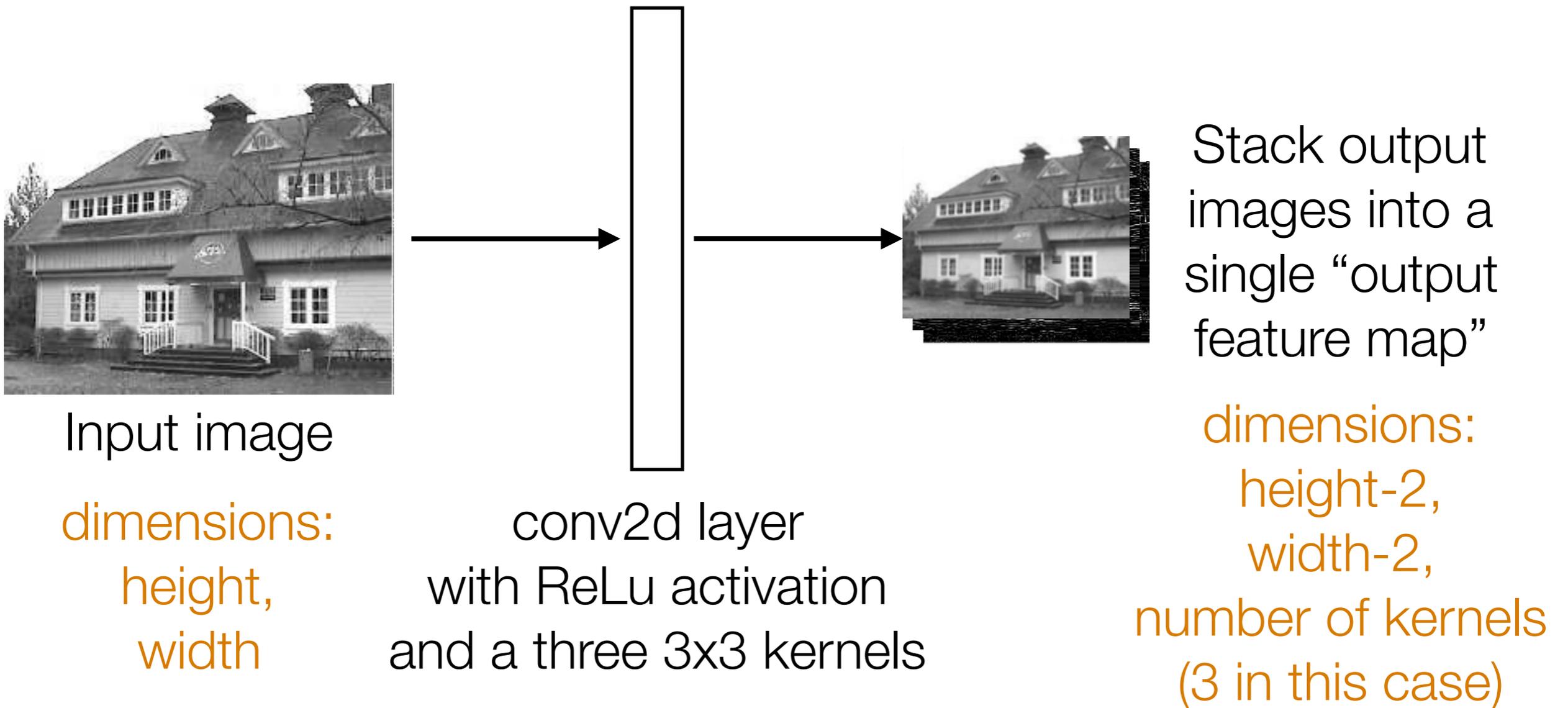
dimensions:
height,
width



conv2d layer
with ReLu activation
and a three 3x3 kernels

Stack output
images into a
single “output
feature map”

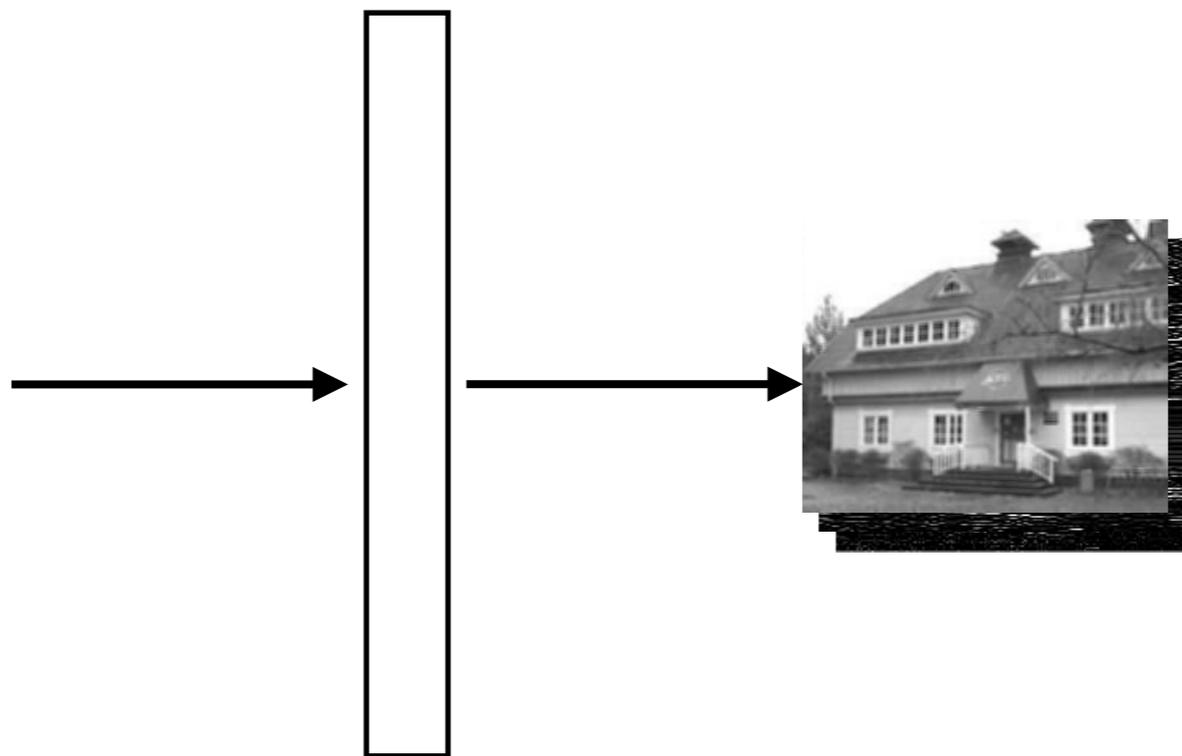
Convolution Layer



Convolution Layer



Input image
dimensions:
height,
width



conv2d layer
with ReLu activation
and k 3x3 kernels



Stack output
images into a
single “output
feature map”

dimensions:
height-2,
width-2,
 k

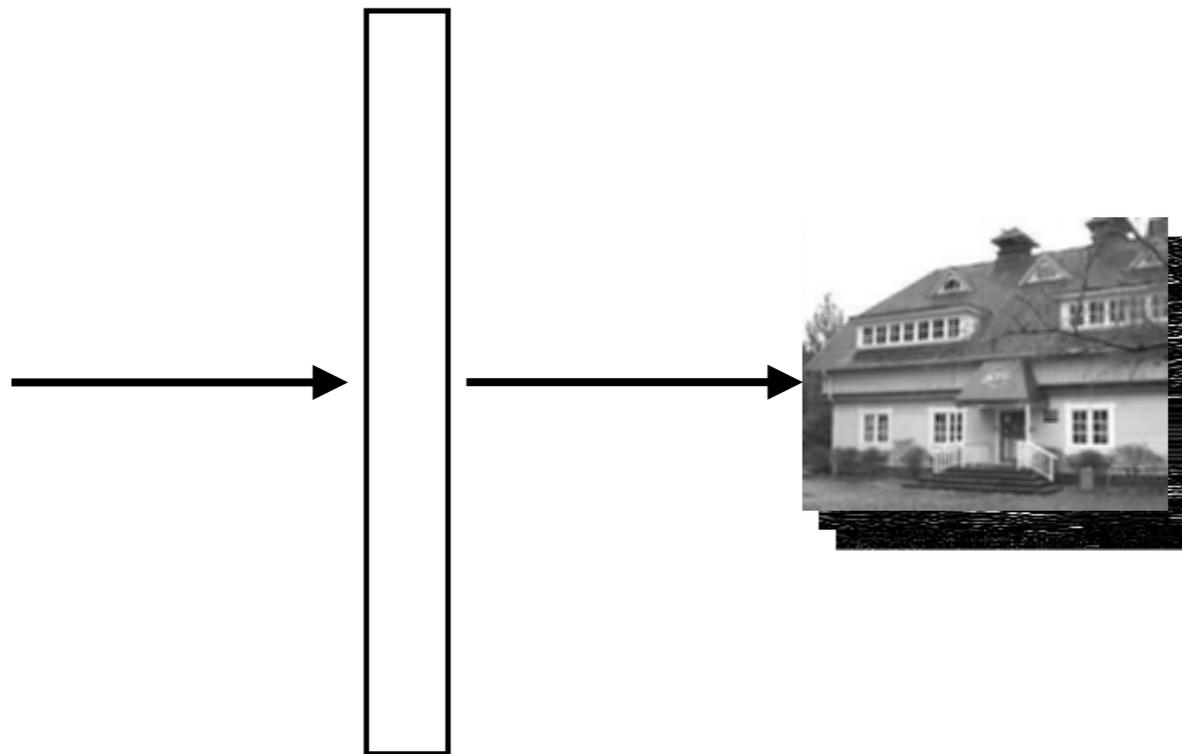
Convolution Layer



Input image

dimensions:
height,
width,

depth d (# channels)



conv2d layer
with ReLu activation
and k $3 \times 3 \times d$ kernels



Stack output
images into a
single “output
feature map”

dimensions:
height-2,
width-2,
 k

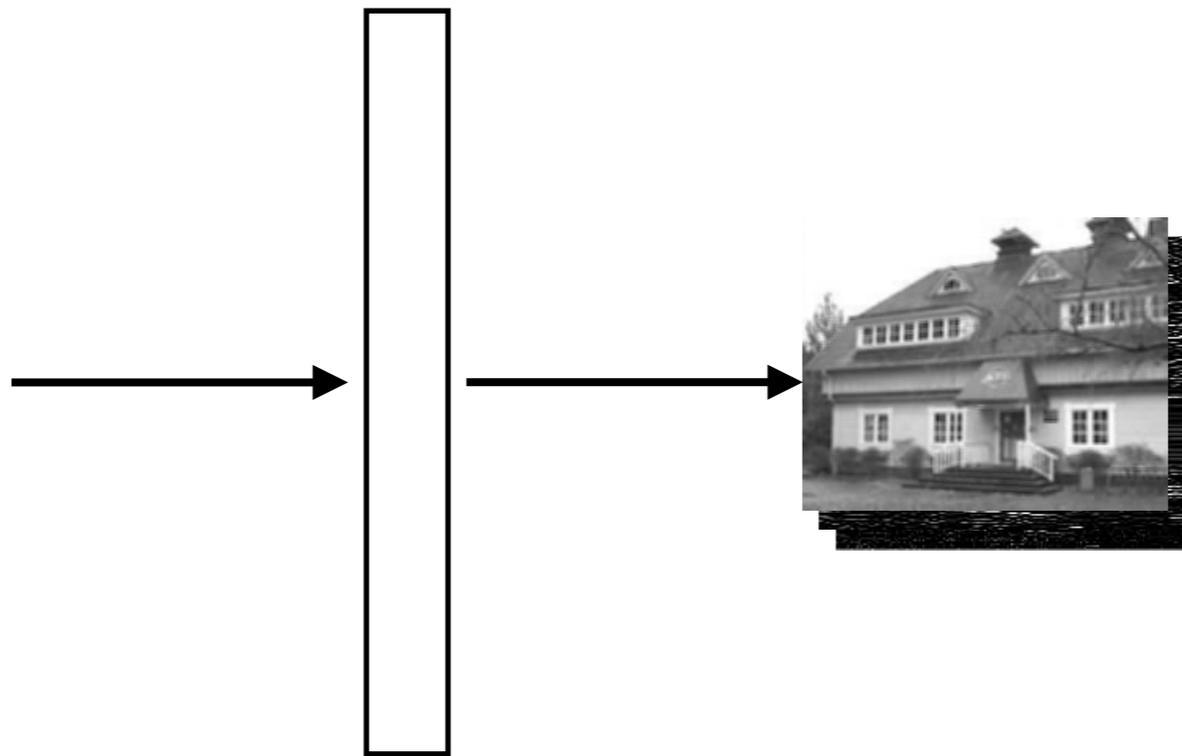
Convolution Layer



Input image

dimensions:
height,
width,

depth d (# channels)



conv2d layer
with ReLu activation
and k $3 \times 3 \times d$ kernels

technical detail: there's
also a bias vector



Stack output
images into a
single “output
feature map”

dimensions:
height-2,
width-2,
 k

Pooling

Pooling

- Aggregate local information

Pooling

- Aggregate local information
- Produces a smaller image
(each resulting pixel captures some “global” information)

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

$$\begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline 2 & 2 & 2 \\ \hline -1 & -1 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 3 & 1 & 0 \\ \hline 1 & 1 & 1 & 3 & 3 \\ \hline 0 & 0 & -2 & -4 & -4 \\ \hline 1 & 1 & 1 & 3 & 3 \\ \hline 0 & 1 & 3 & 1 & 0 \\ \hline \end{array}$$

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

1	

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

1	3

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

1	3
1	

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

1	3
1	3

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image
after ReLU

1	3
1	3

Output after
max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?



1	3
1	3

Output after max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?



1	3
1	3

Output after max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?



1	3
1	3

Output after max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?



1	3
1	3

Output after max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?

In this example: 1 pixel in max pooling output captures information from 16 input pixels!



1	3
1	3

Output after max pooling

Max Pooling

0	0	0	0	0	0	0
0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	1	0	0	0
0	1	1	1	1	1	0
0	0	1	1	1	0	0
0	0	0	0	0	0	0

Input image

-1	-1	-1
2	2	2
-1	-1	-1

*

=

0	1	3	1	0
1	1	1	3	3
0	0	-2	-4	-4
1	1	1	3	3
0	1	3	1	0

0	1	3	1	0
1	1	1	3	3
0	0	0	0	0
1	1	1	3	3
0	1	3	1	0

Output image after ReLU

What numbers were involved in computing this 1?

In this example: 1 pixel in max pooling output captures information from 16 input pixels!

Example: applying max pooling again results in a single pixel that captures info from entire input image!



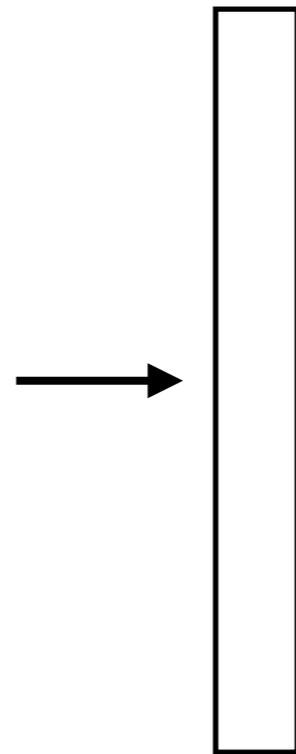
1	3
1	3

Output after max pooling

Basic Building Block of CNN's

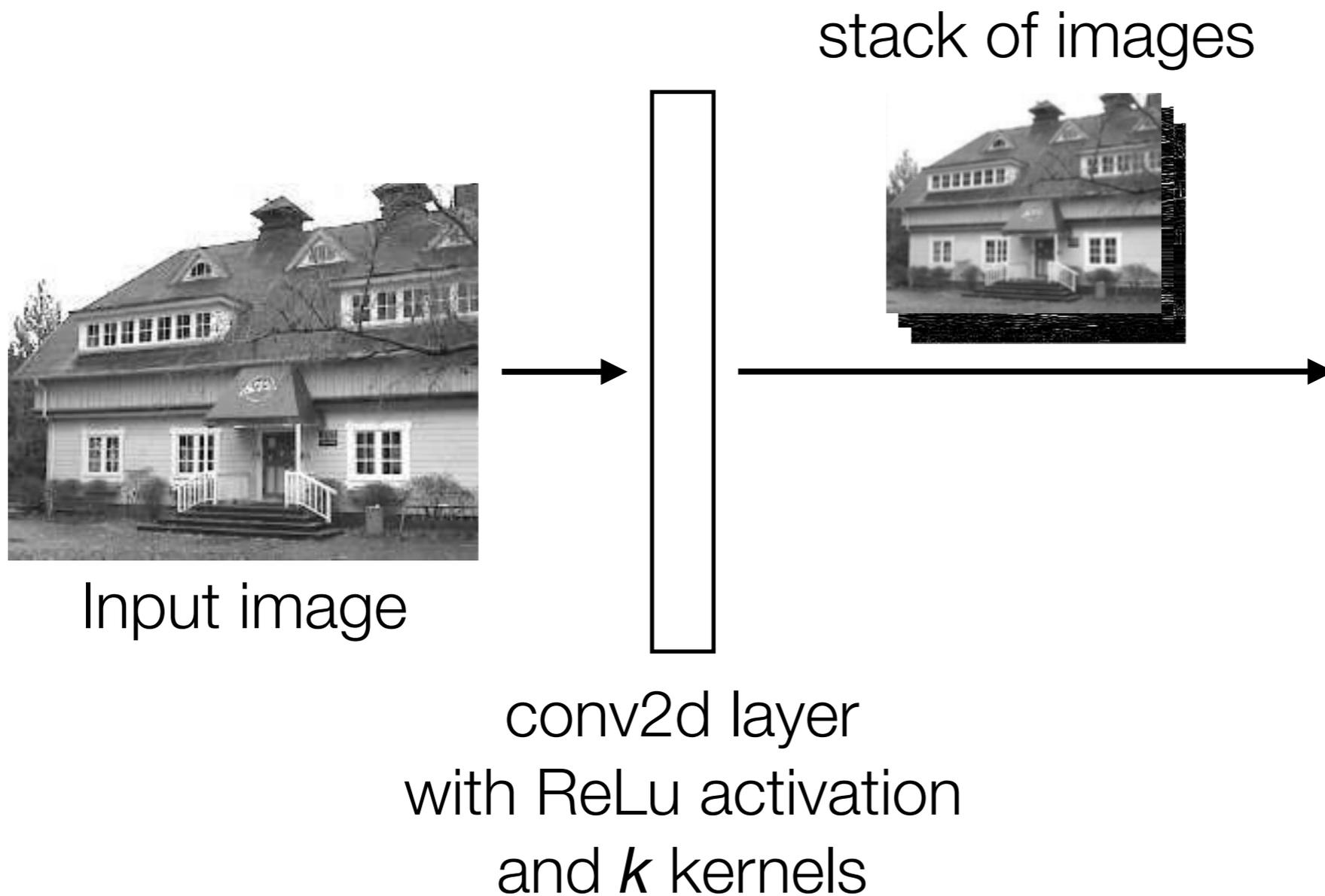


Input image

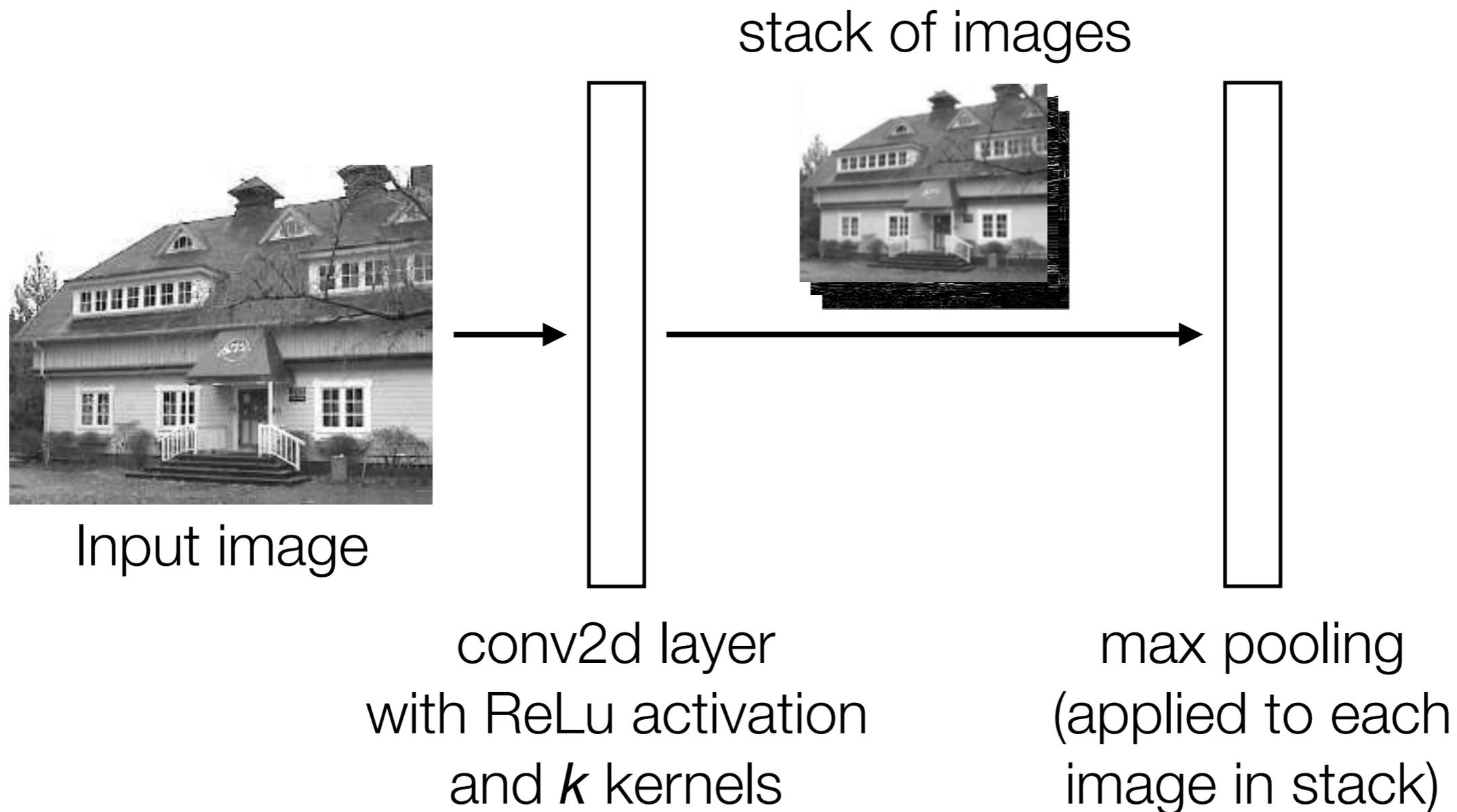


conv2d layer
with ReLu activation
and k kernels

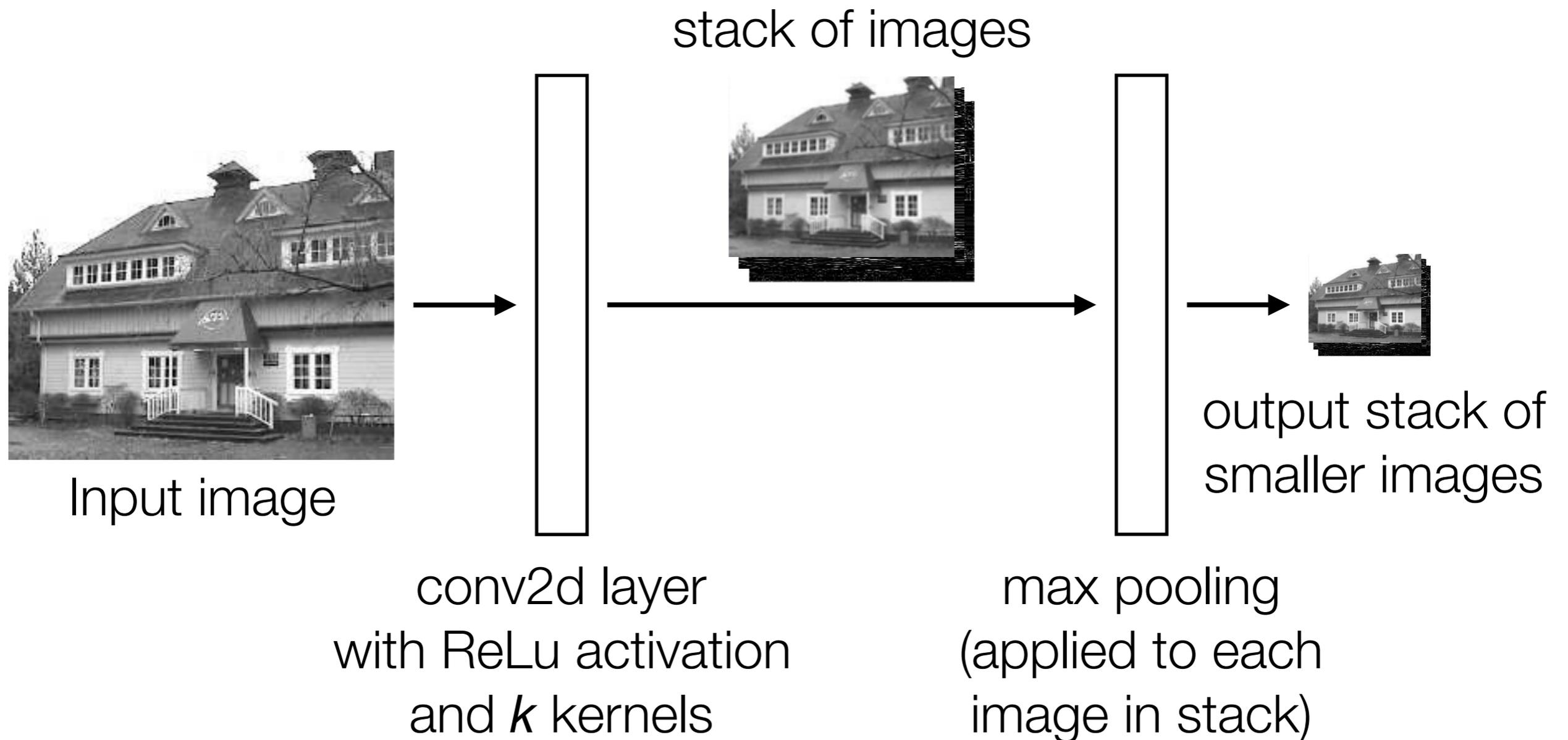
Basic Building Block of CNN's



Basic Building Block of CNN's

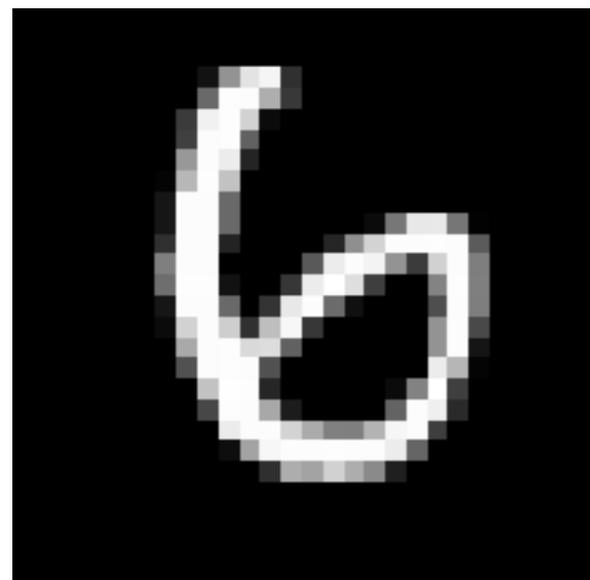


Basic Building Block of CNN's



Handwritten Digit Recognition

Training label: 6



28x28 image

length 784 vector
(784 input neurons)

Learning this neural net means learning parameters of both dense layers!



dense layer with 512 neurons, ReLU activation

dense layer with 10 neurons, softmax activation

Loss/"error"

Popular loss function for classification (> 2 classes): **categorical cross entropy**

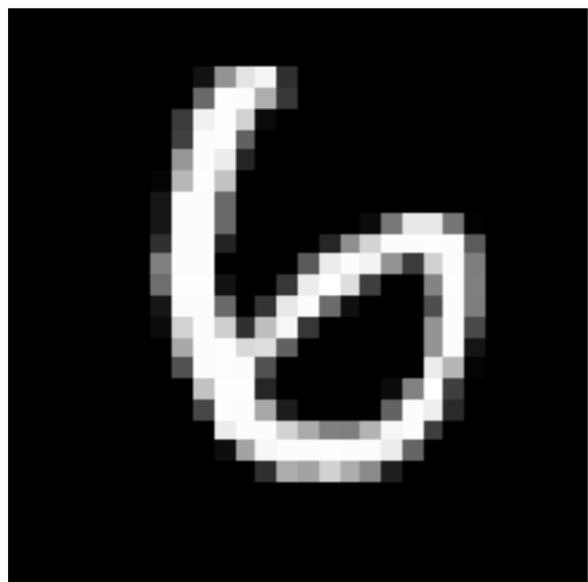
$$\log \frac{1}{\text{Pr}(\text{digit } 6)}$$

Error is averaged across training examples

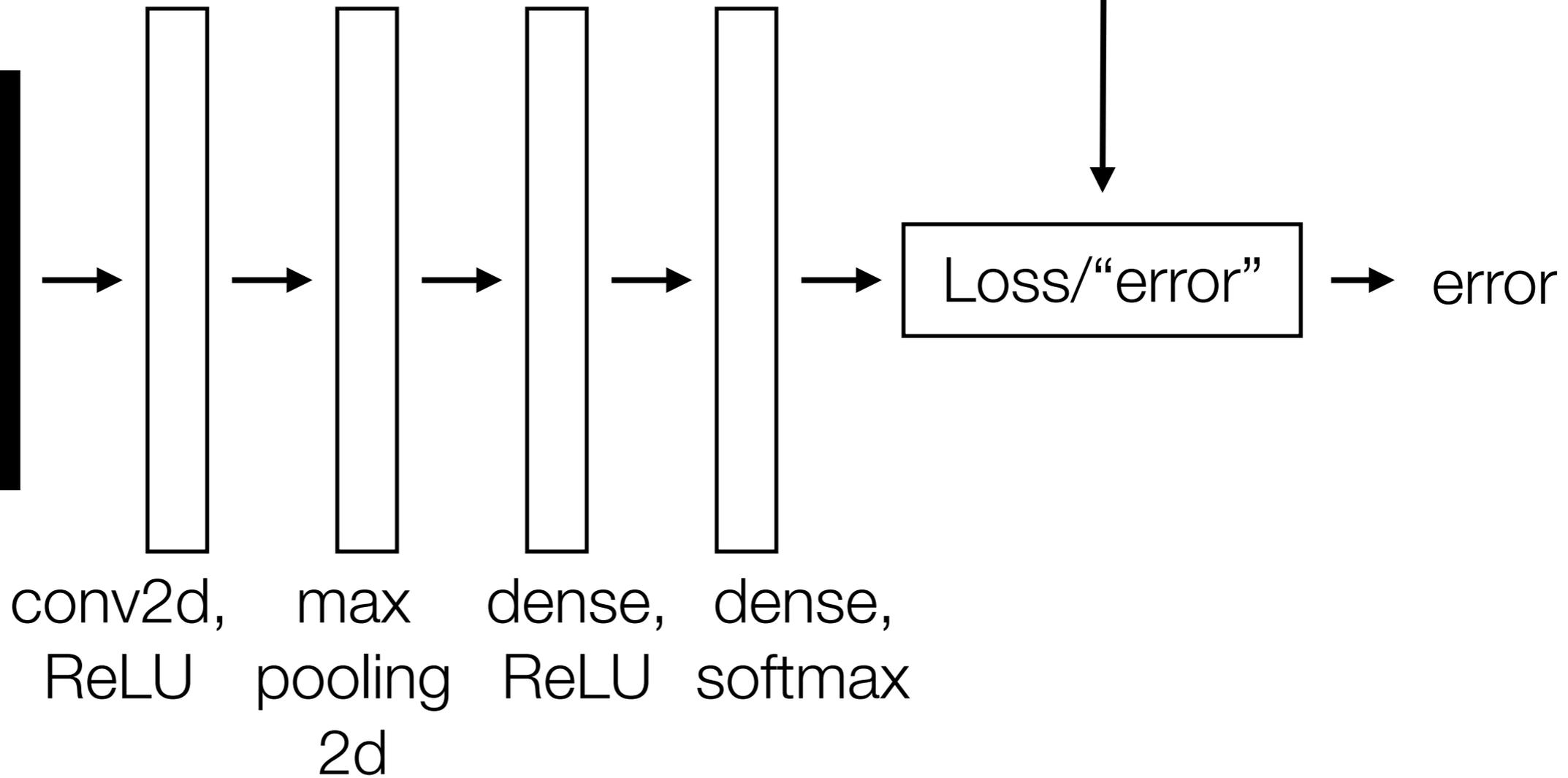
error

Handwritten Digit Recognition

Training label: 6

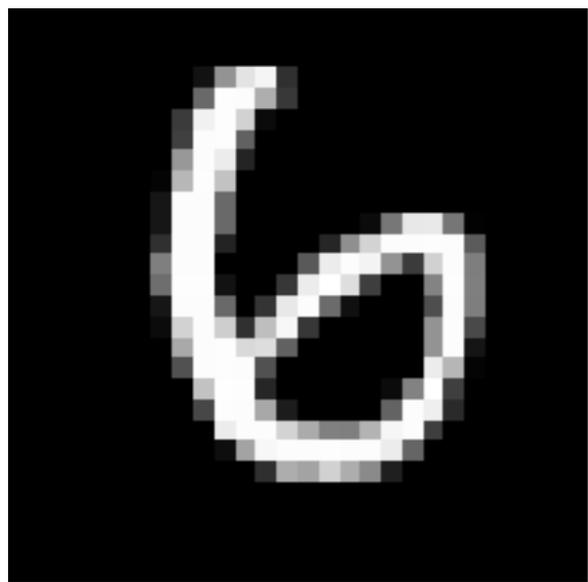


28x28 image

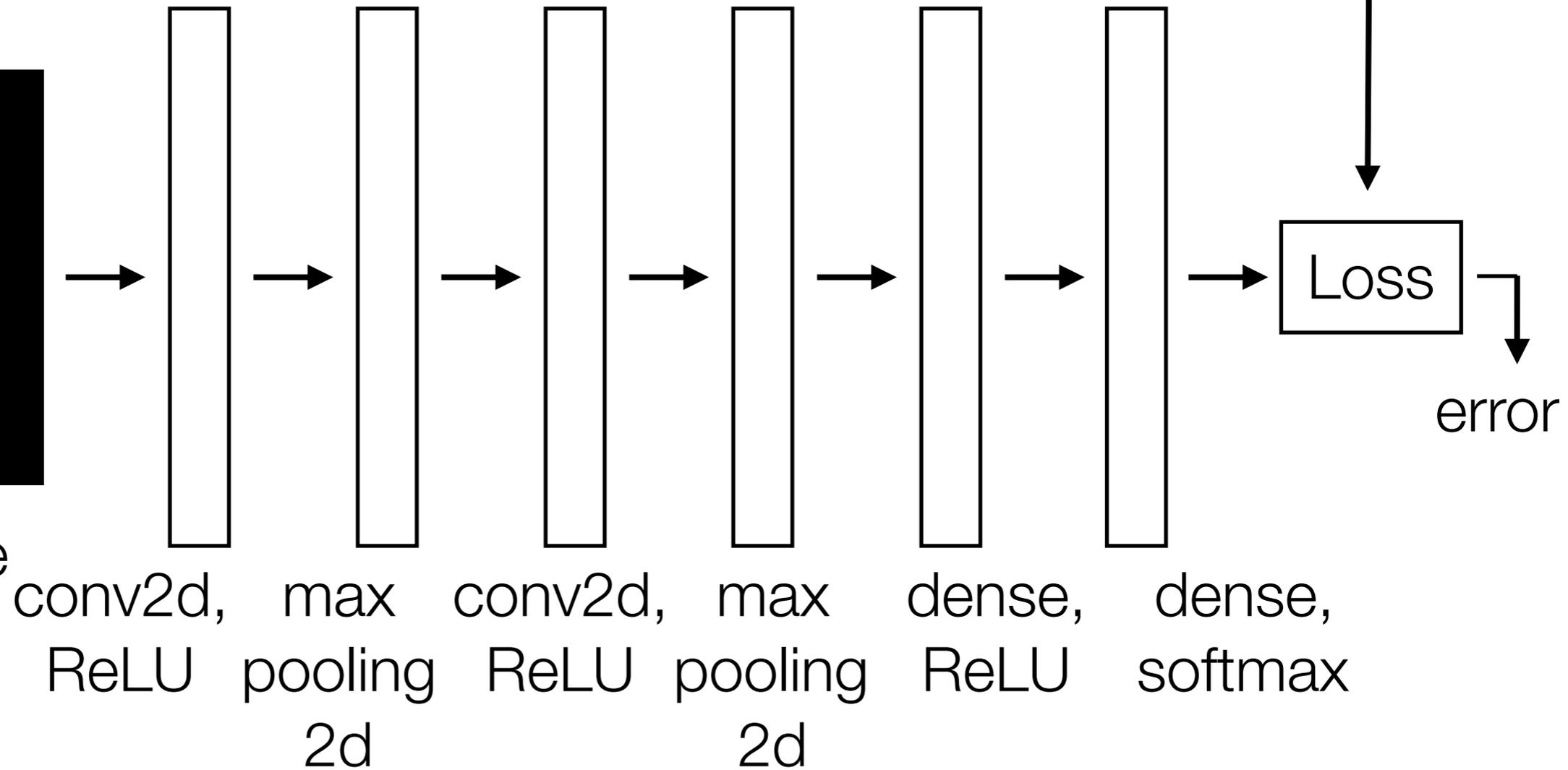


Handwritten Digit Recognition

Training label: 6



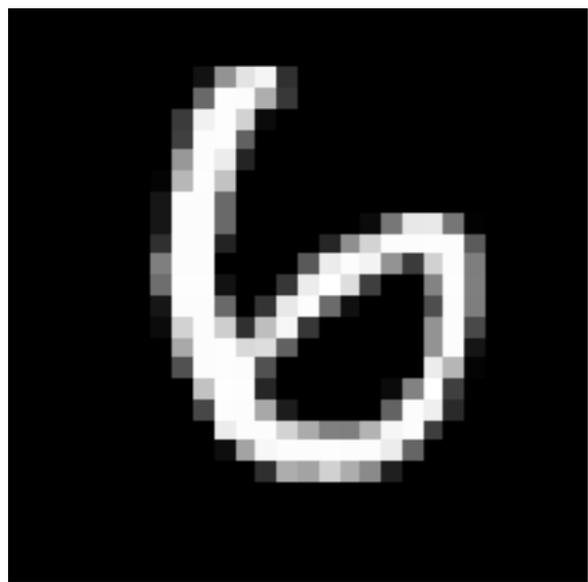
28x28 image



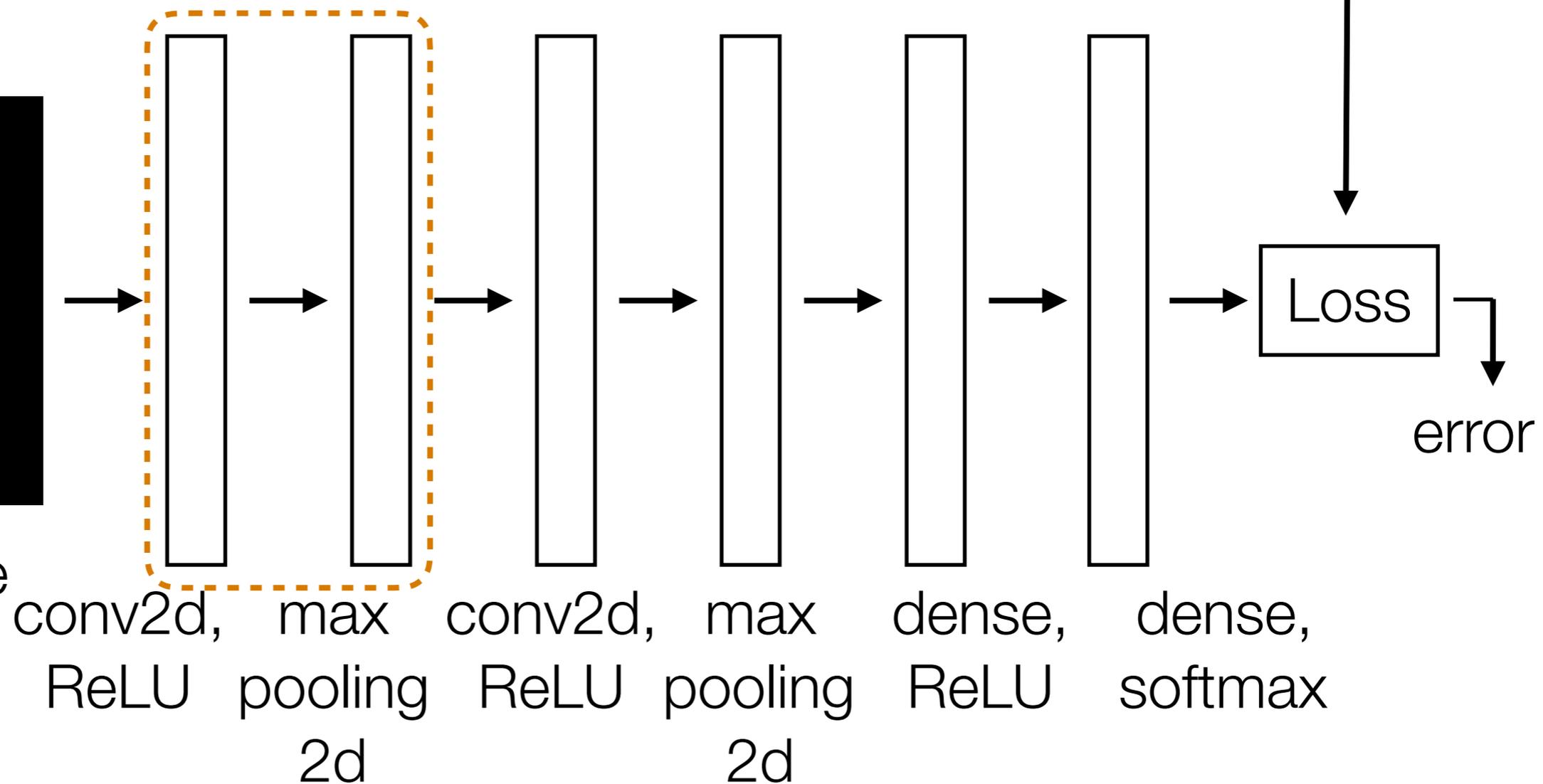
Handwritten Digit Recognition

Training label: 6

extract low-level visual features & aggregate



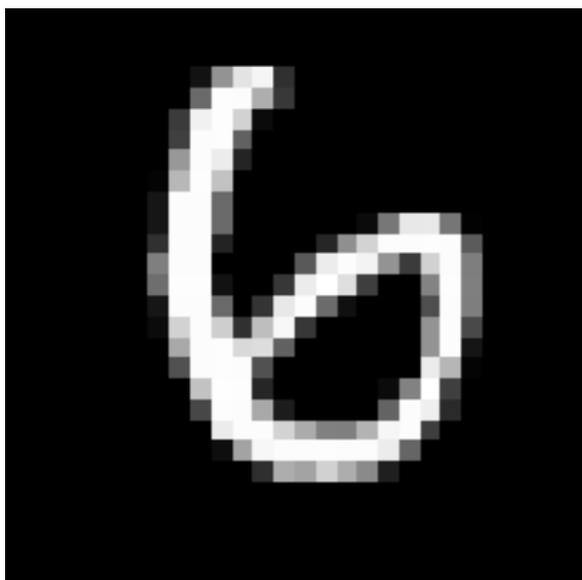
28x28 image



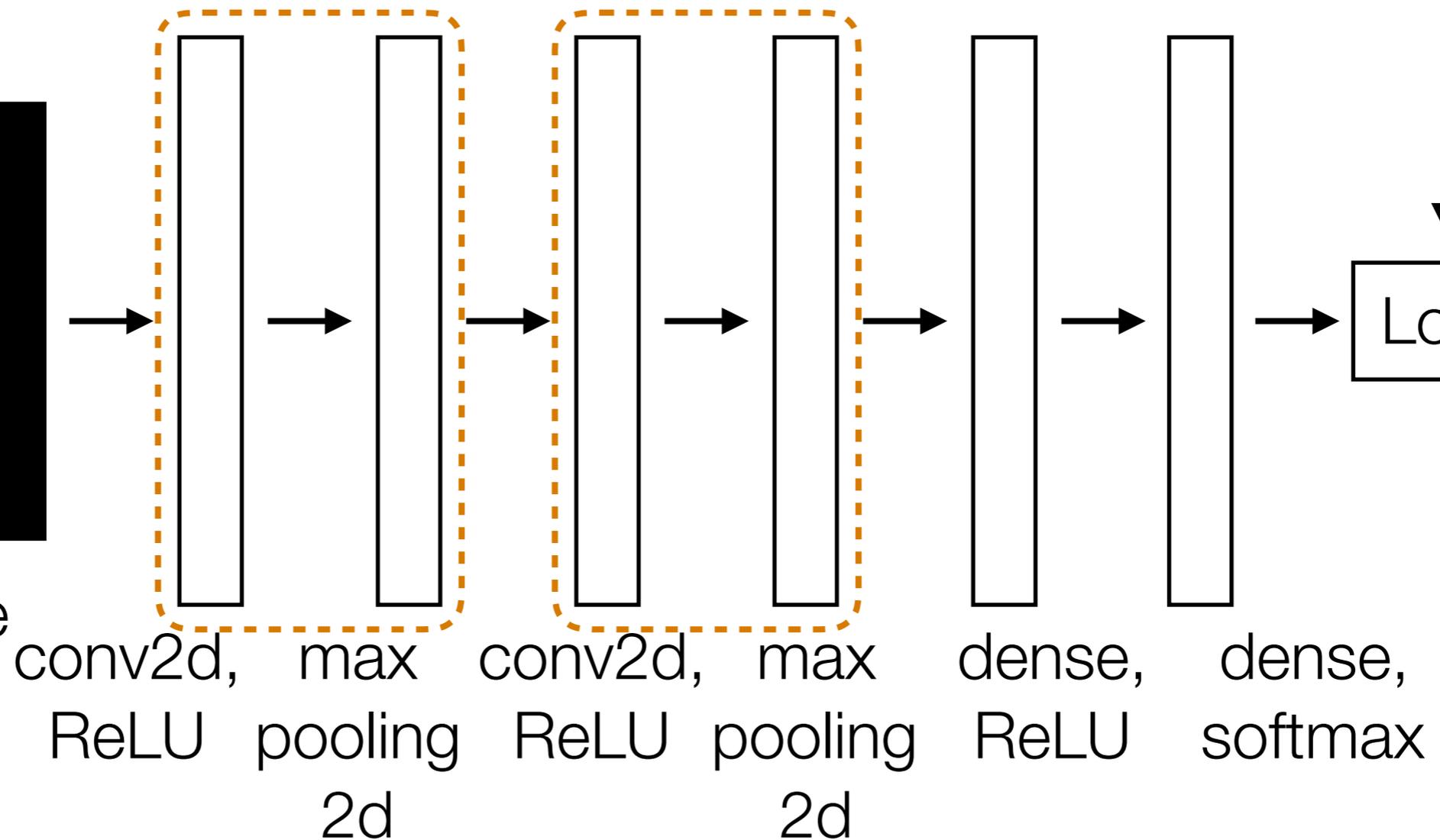
Handwritten Digit Recognition

Training label: 6

extract low-level visual features & aggregate



28x28 image



extract higher-level visual features & aggregate

Loss

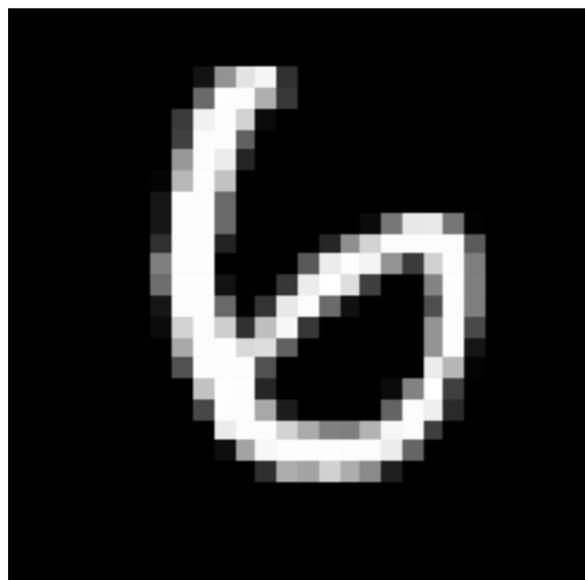
error

Handwritten Digit Recognition

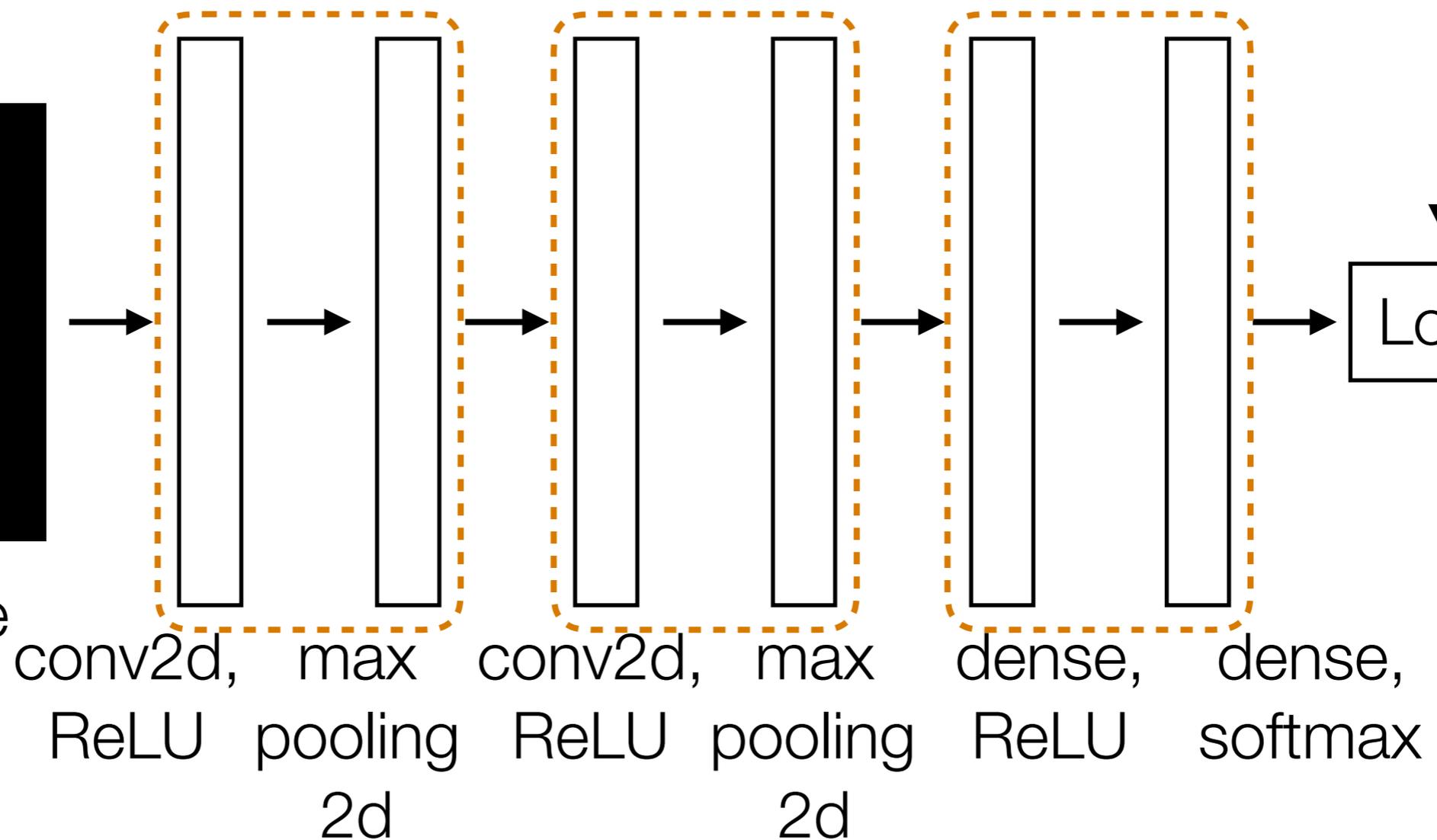
Training label: 6

extract low-level visual features & aggregate

non-vision-specific classification neural net



28x28 image



extract higher-level visual features & aggregate

CNN Demo

CNN's

CNN's

- Learn convolution filters for extracting simple features

CNN's

- Learn convolution filters for extracting simple features
- Max pooling aggregates local information

CNN's

- Learn convolution filters for extracting simple features
- Max pooling aggregates local information
- Can then repeat the above two layers to learn features from increasingly higher-level representations

CNN's

- Learn convolution filters for extracting simple features
- Max pooling aggregates local information
- Can then repeat the above two layers to learn features from increasingly higher-level representations
- Convolution filters are shift-invariant

CNN's

- Learn convolution filters for extracting simple features
- Max pooling aggregates local information
- Can then repeat the above two layers to learn features from increasingly higher-level representations
- Convolution filters are shift-invariant
- In terms of invariance to an object shifting within the input image, this is roughly achieved by pooling

Recurrent Neural Networks (RNNs)

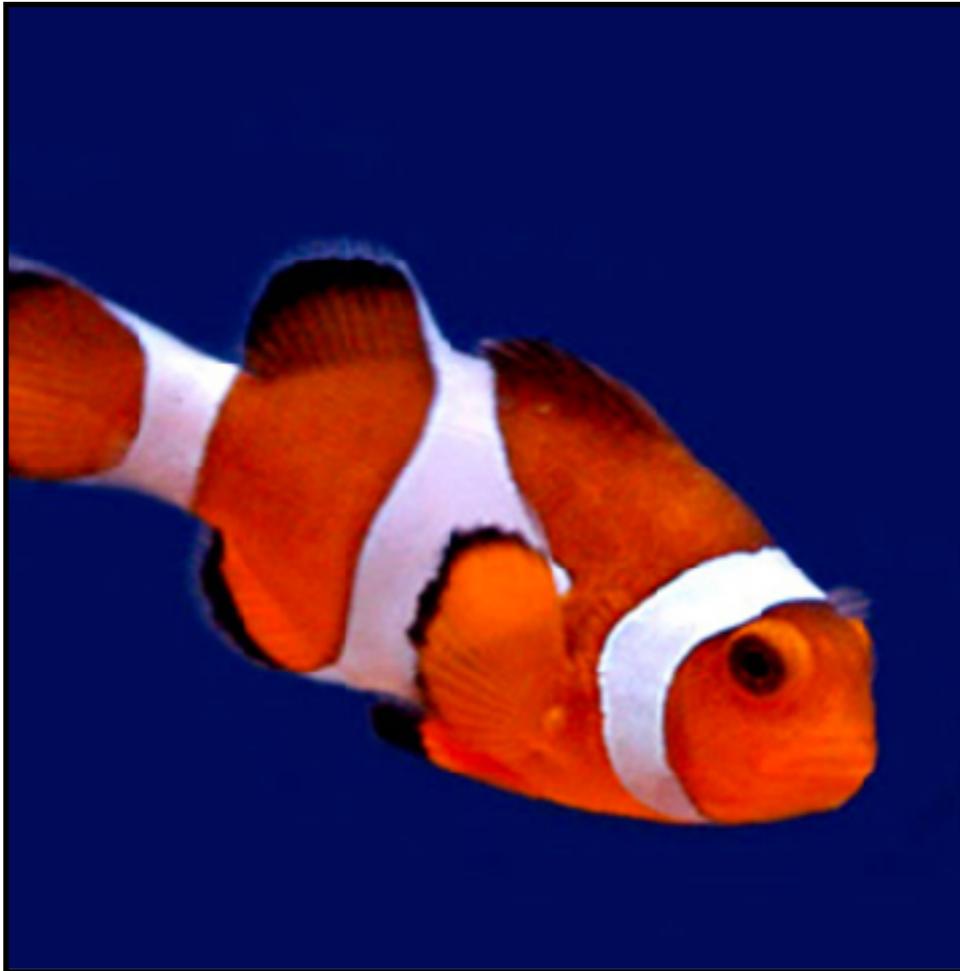
RNNs

RNNs

What we've seen so far are "feedforward" NNs

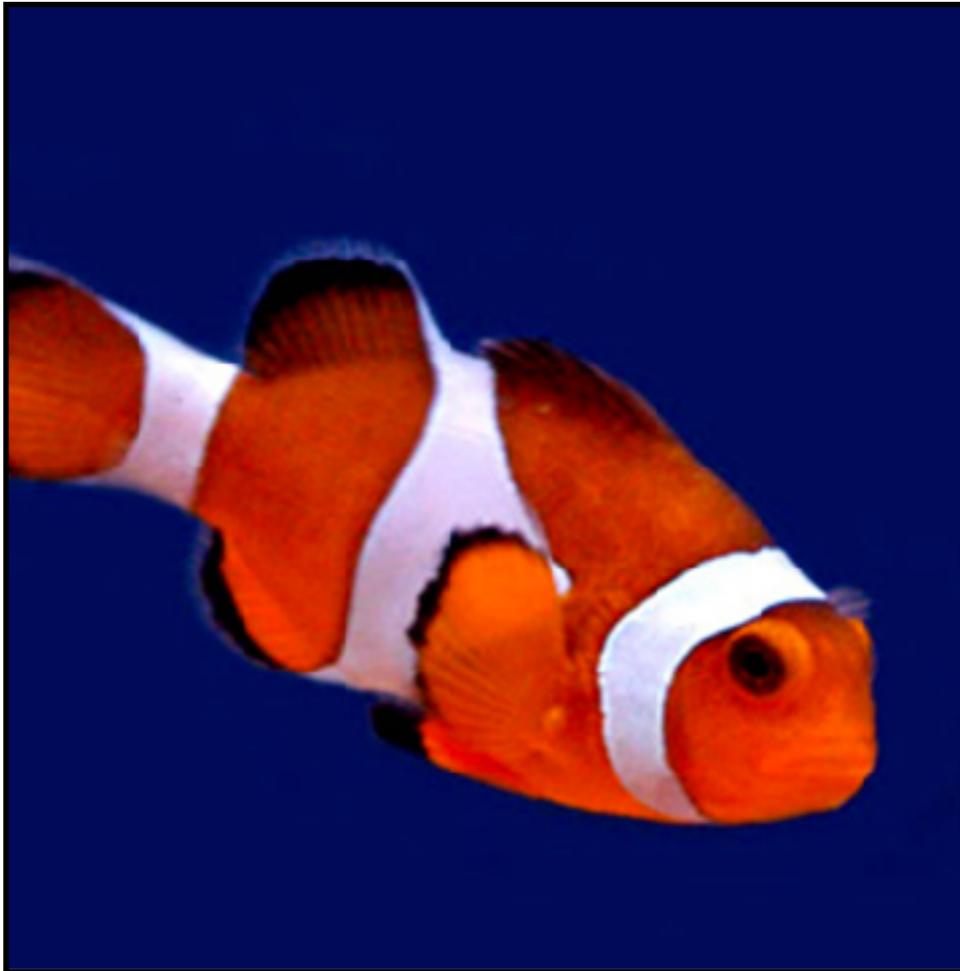
RNNs

What we've seen so far are "feedforward" NNs



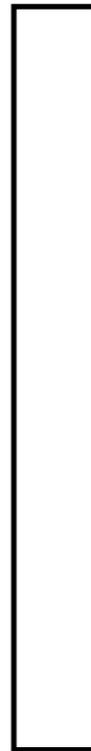
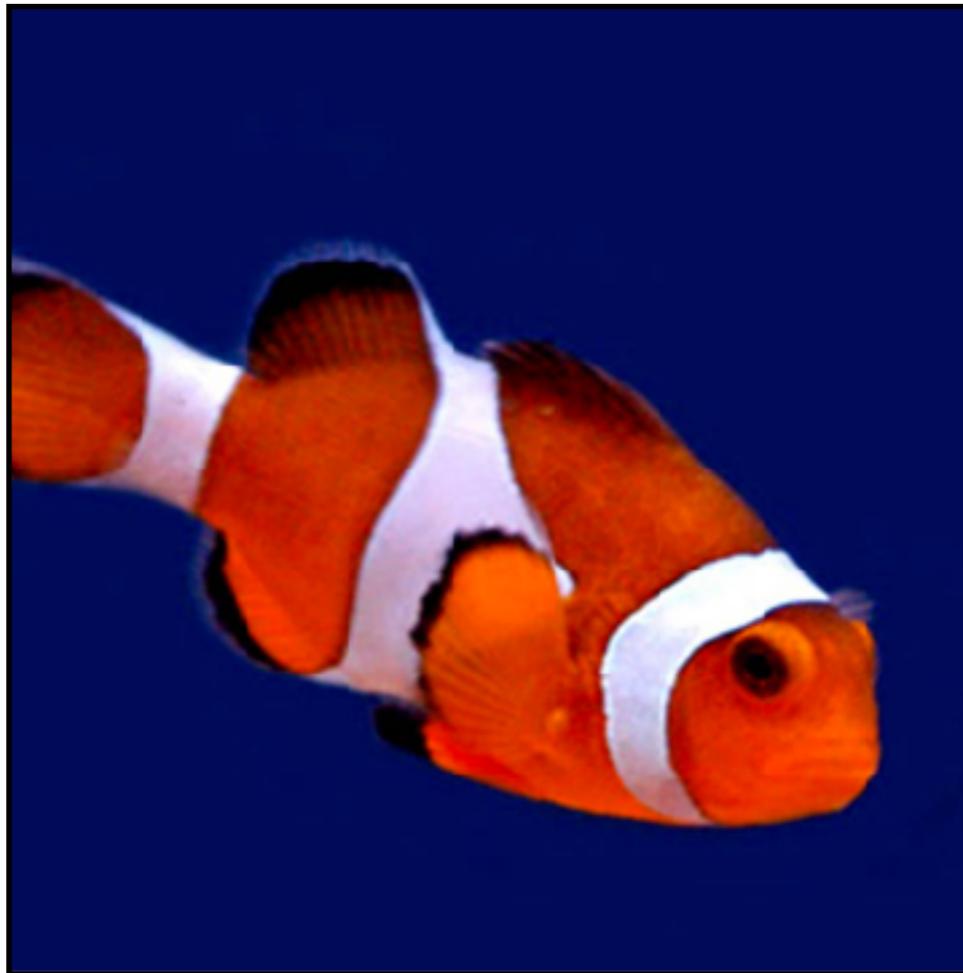
RNNs

What we've seen so far are "feedforward" NNs



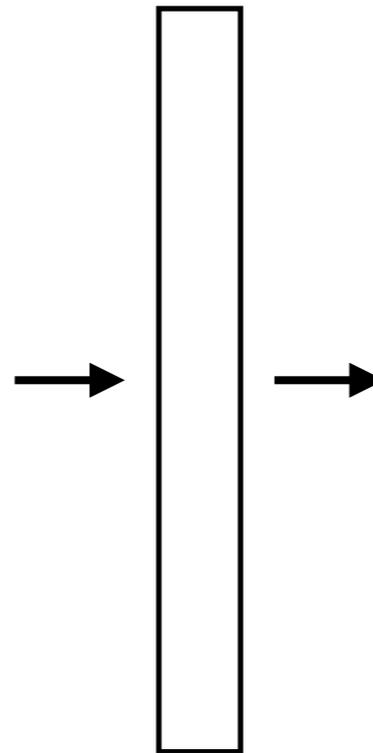
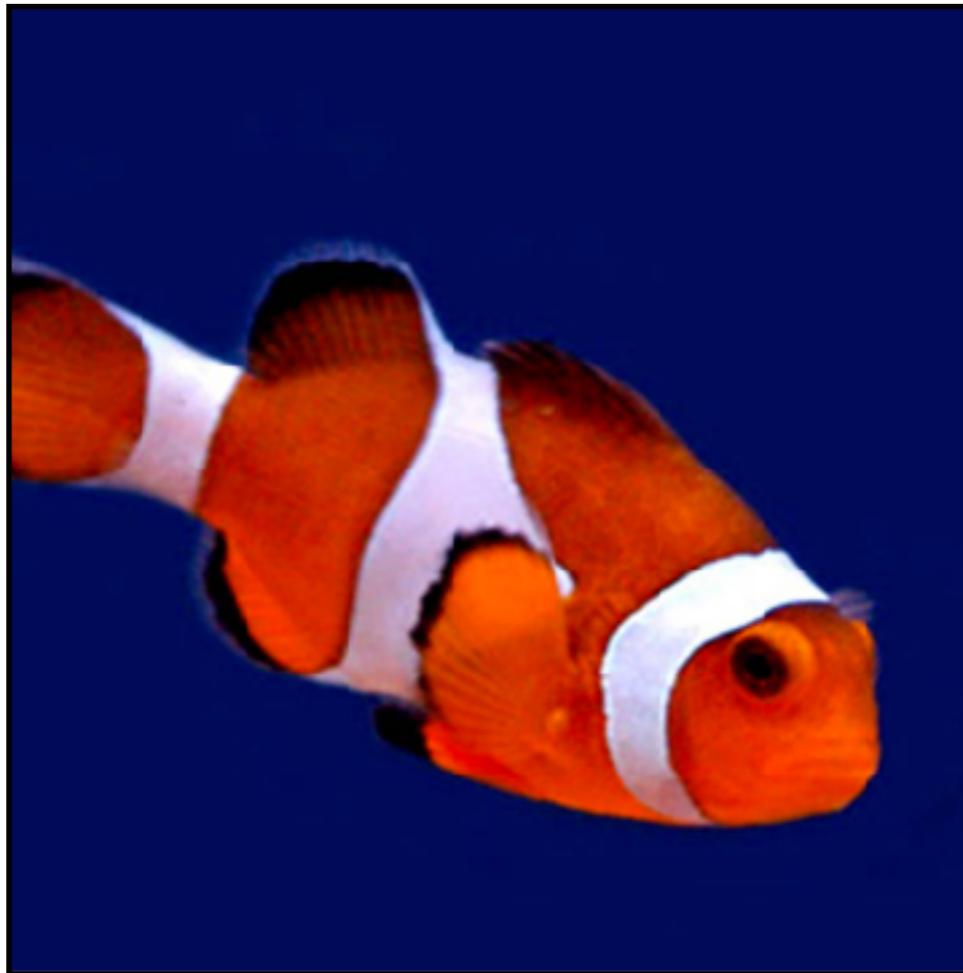
RNNs

What we've seen so far are "feedforward" NNs



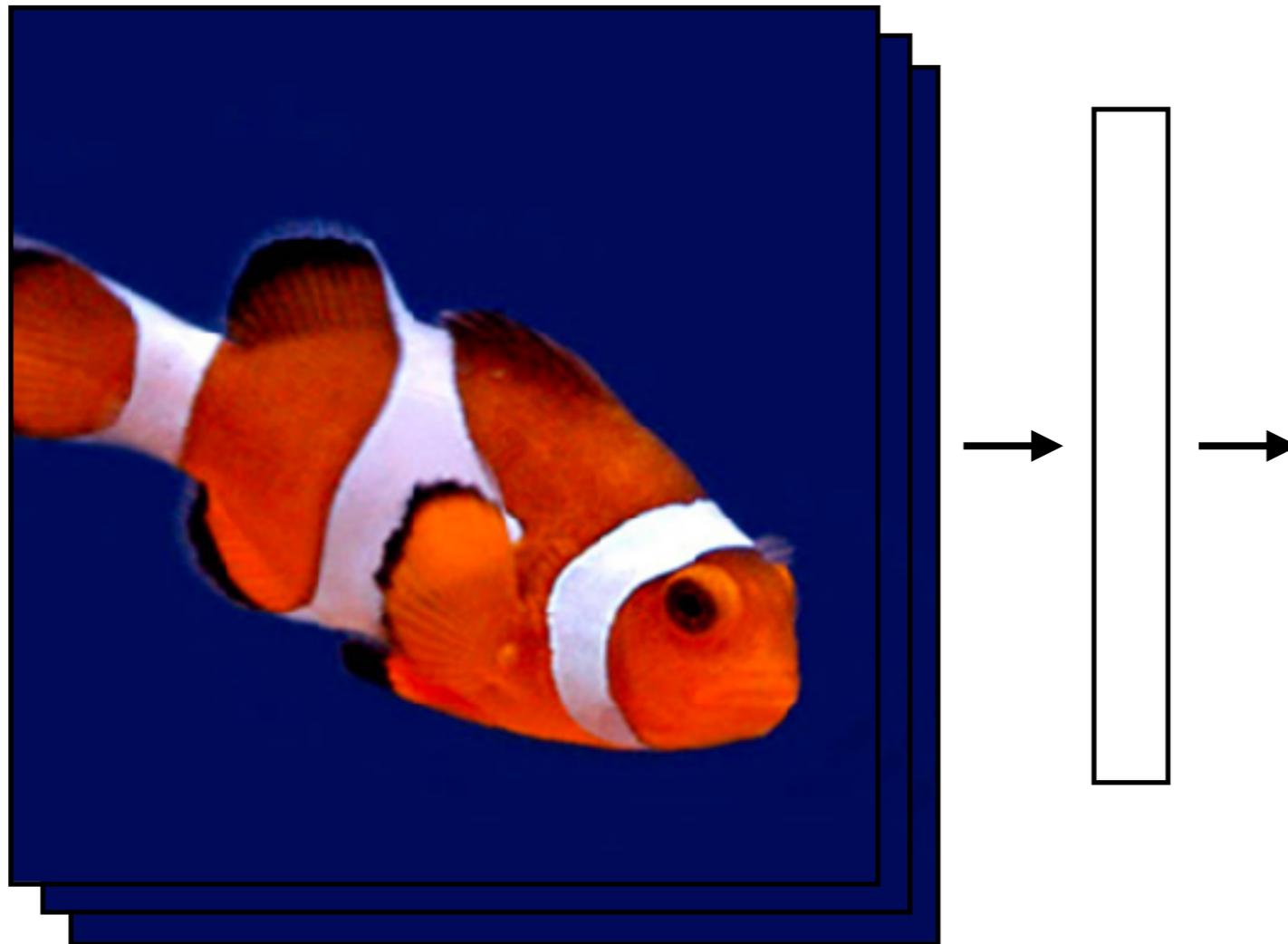
RNNs

What we've seen so far are "feedforward" NNs



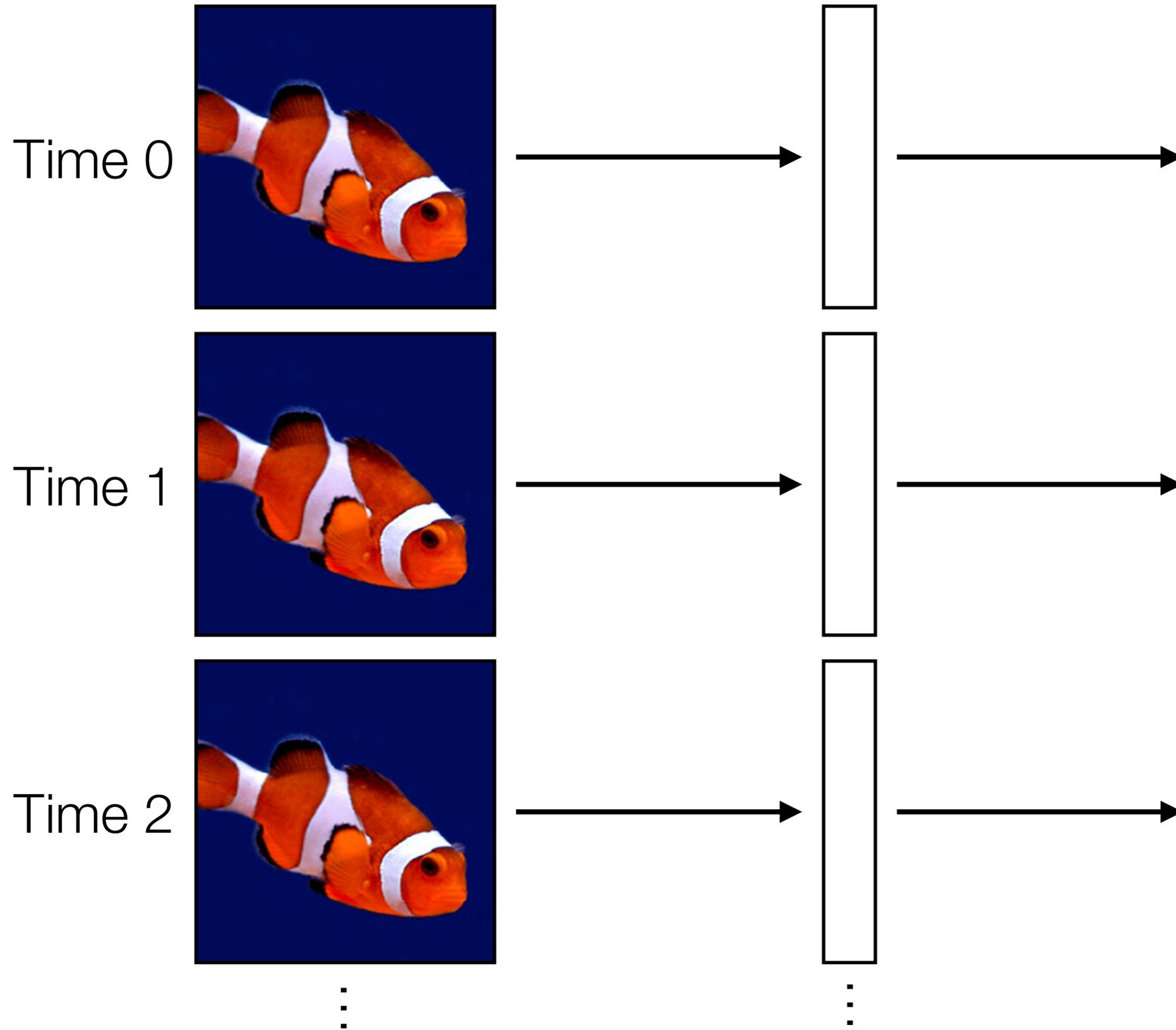
RNNs

What we've seen so far are "feedforward" NNs



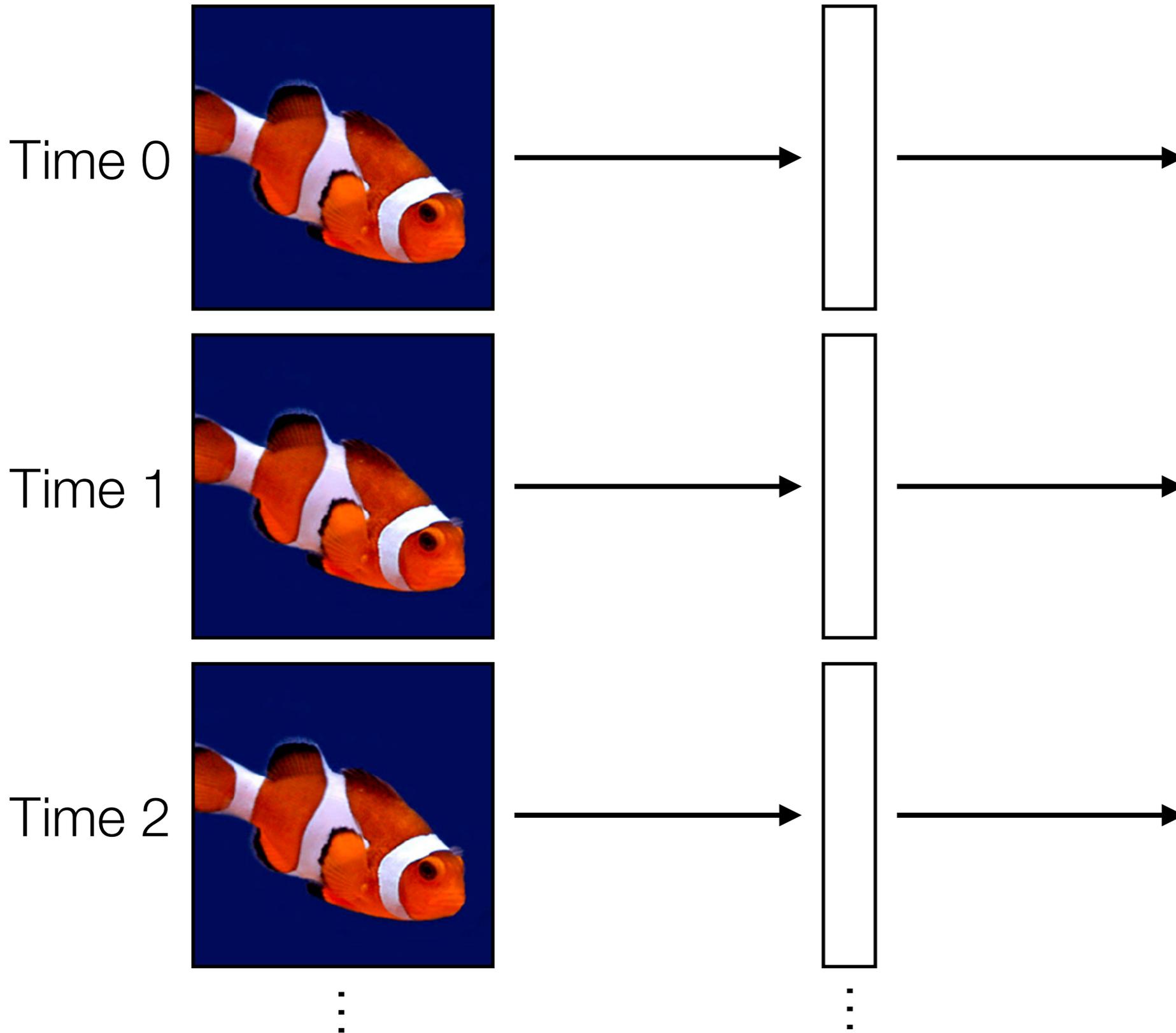
What if we had a video?

RNNs



RNNs

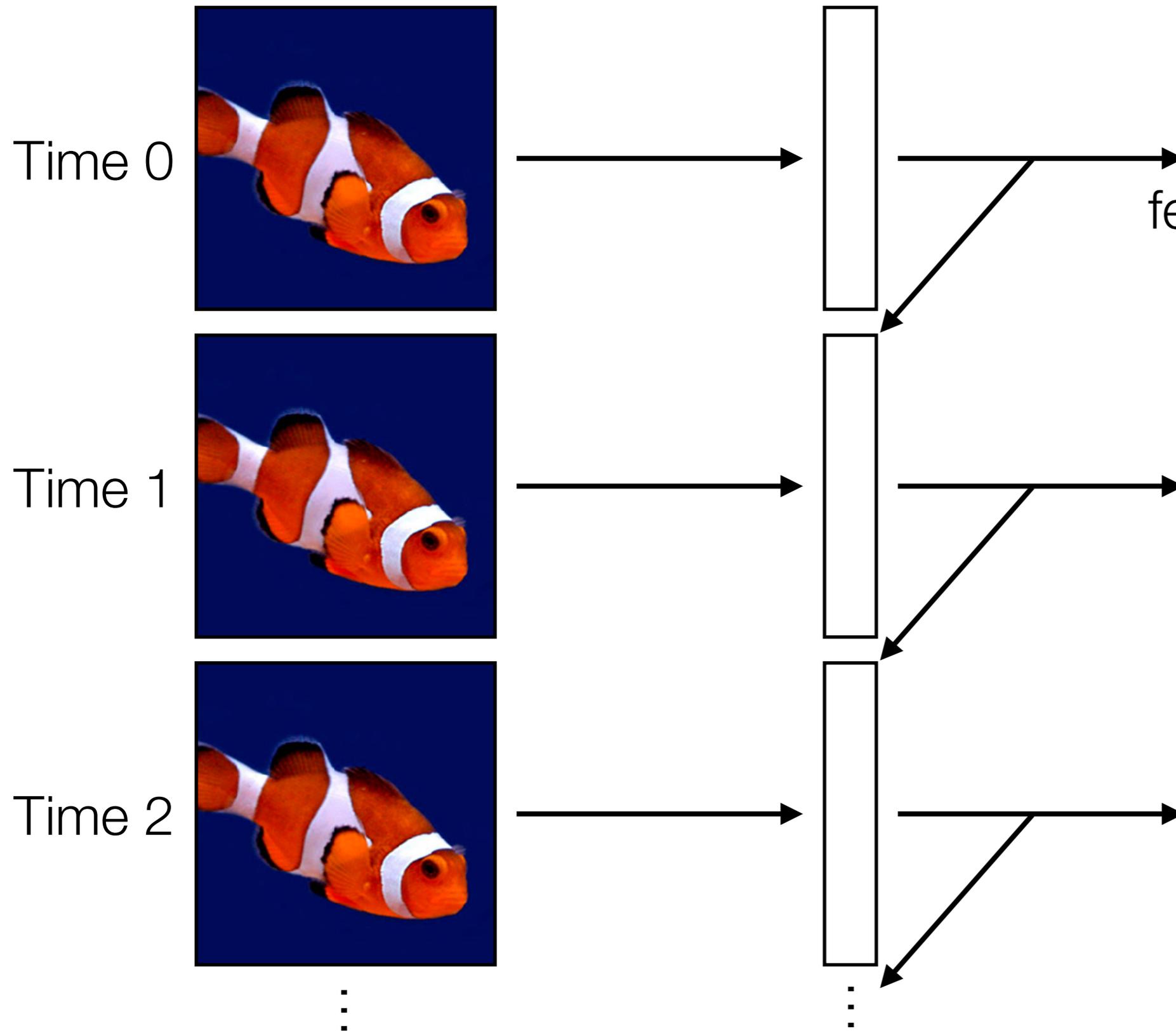
Feedforward NN's:
treat each video frame
separately



RNNs

Feedforward NN's:
treat each video frame
separately

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



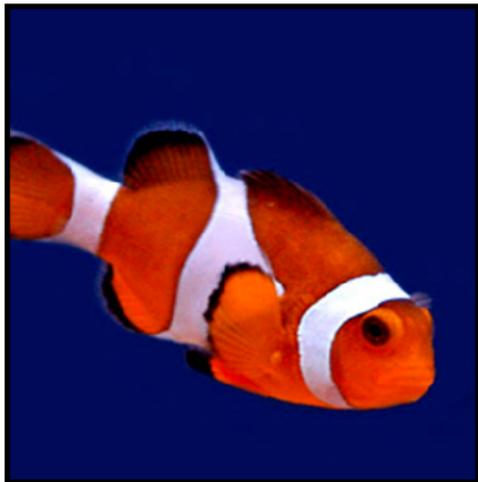
RNNs

Feedforward NN's:
treat each video frame
separately

Time 0



Time 1



Time 2



⋮

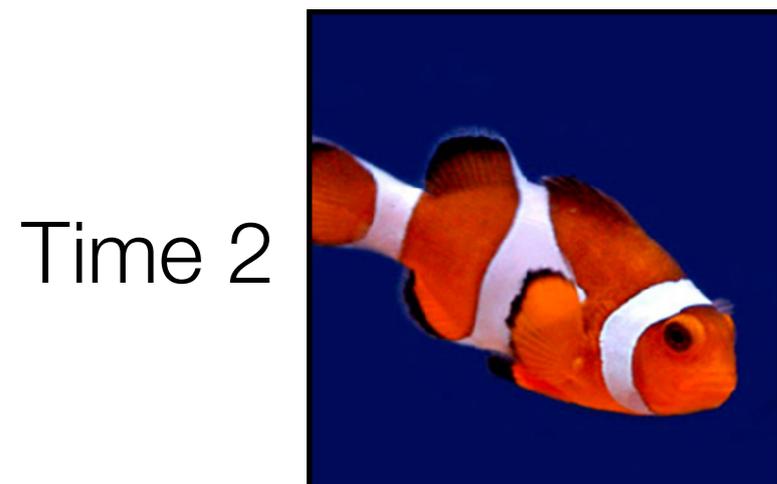
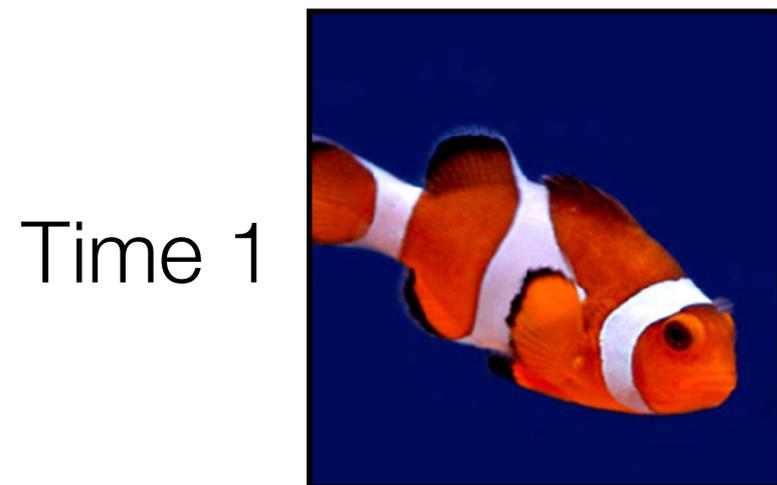
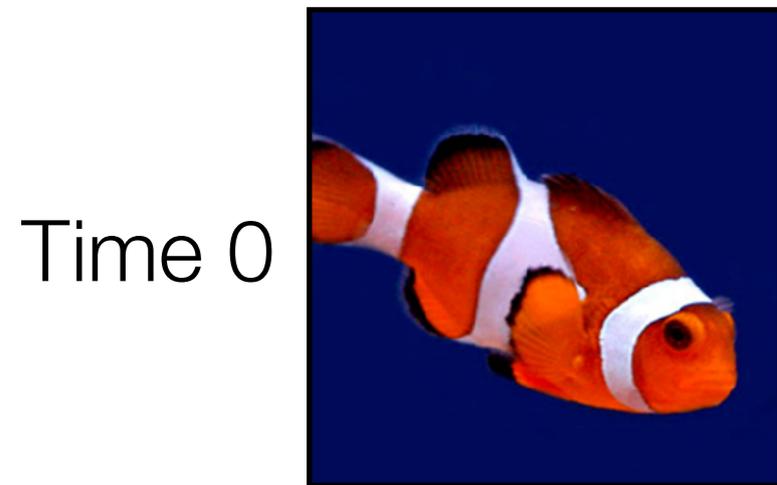
⋮

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step

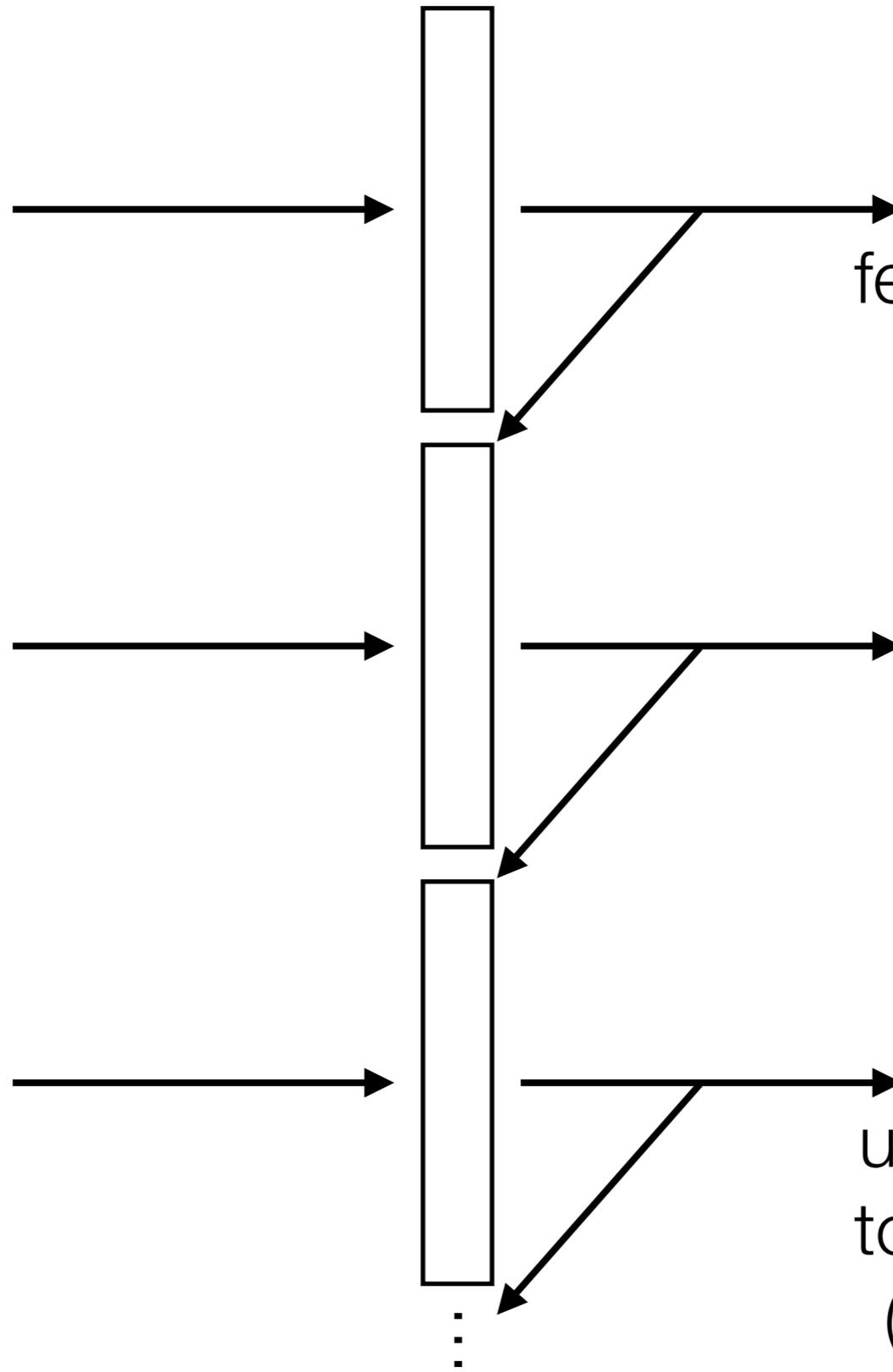
In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

RNNs

Feedforward NN's:
treat each video frame
separately



⋮



RNN's:
feed output at previous
time step as input to
RNN layer at current
time step

In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

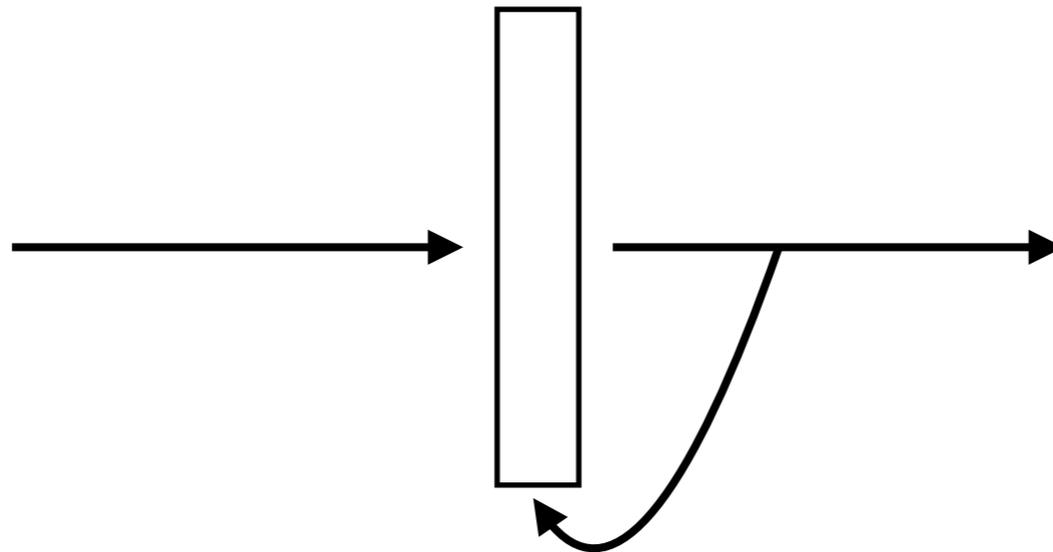
RNNs

Feedforward NN's:
treat each video frame
separately

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



Time series



In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

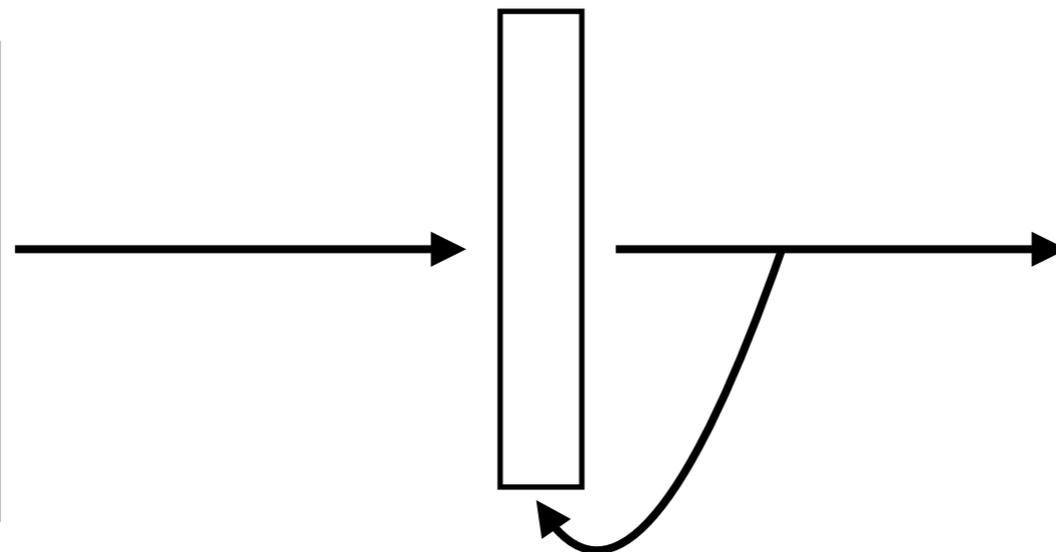
RNNs

Feedforward NN's:
treat each video frame
separately

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



Time series



LSTM layer

In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

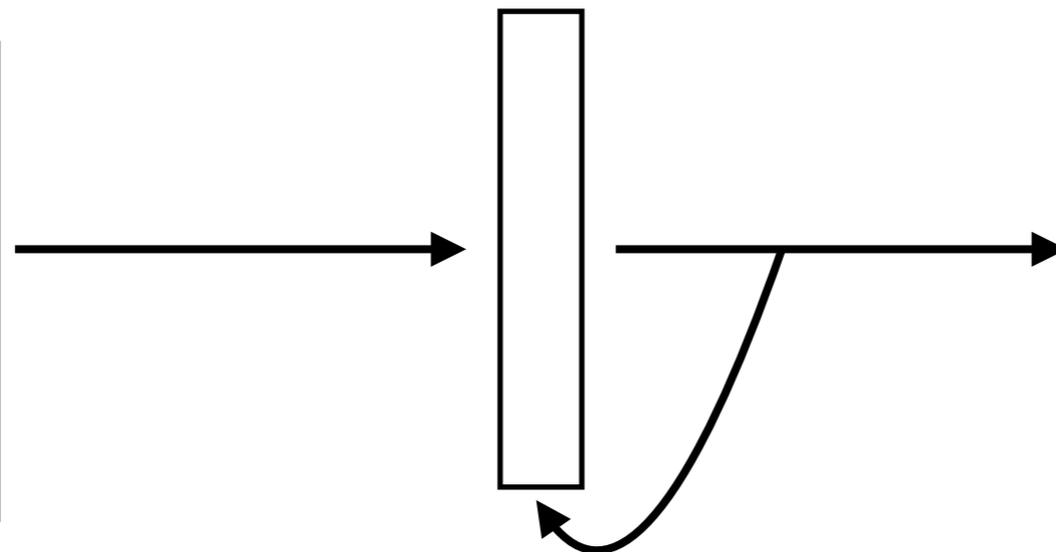
RNNs

Feedforward NN's:
treat each video frame
separately

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



Time series



LSTM layer

like a dense layer
that has memory

In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

RNNs

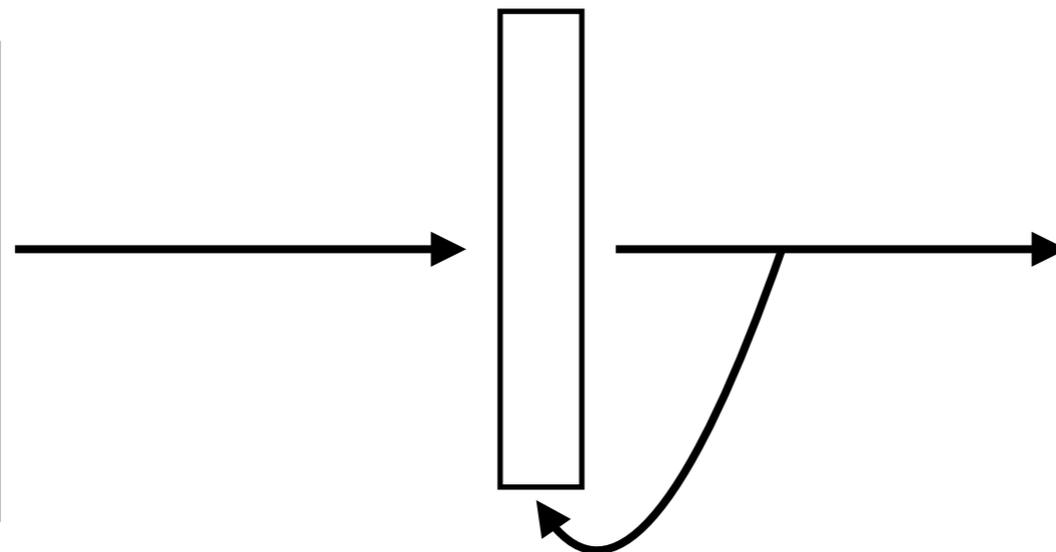
Feedforward NN's:
treat each video frame
separately

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step

readily chains together with
other neural net layers



Time series



LSTM layer

like a dense layer
that has memory

In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

RNNs

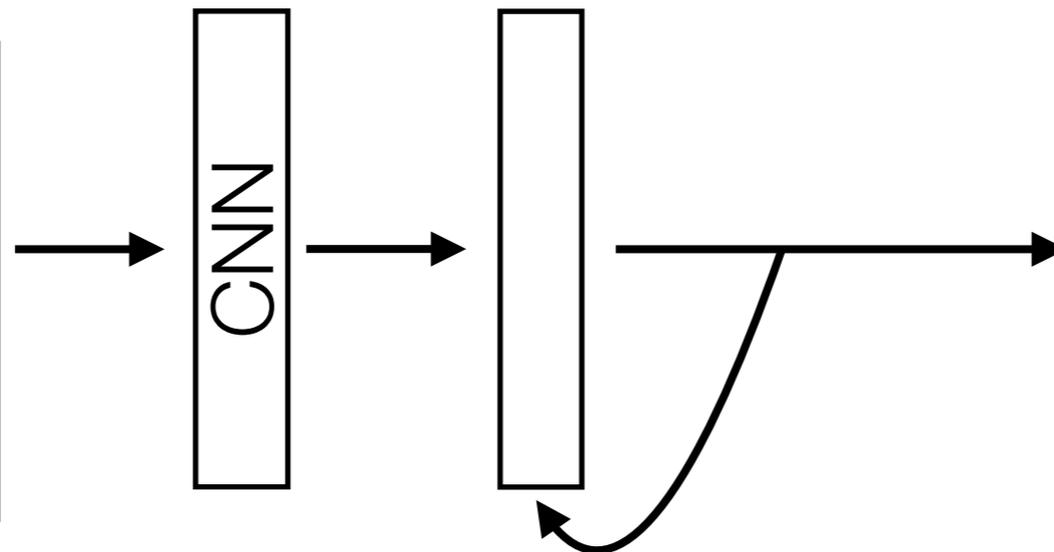
Feedforward NN's:
treat each video frame
separately

readily chains together with
other neural net layers

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



Time series



LSTM layer

like a dense layer
that has memory

In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

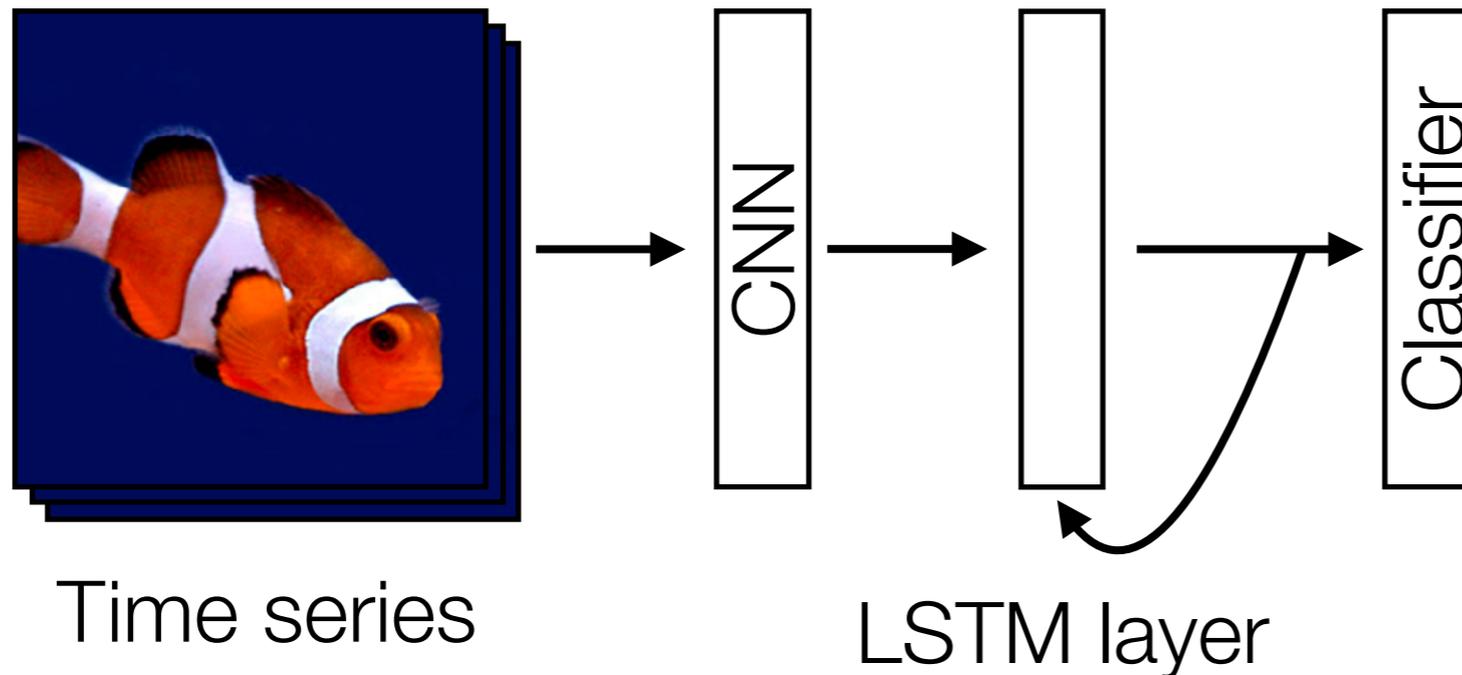
Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

RNNs

Feedforward NN's:
treat each video frame
separately

readily chains together with
other neural net layers

RNN's:
feed output at previous
time step as input to
RNN layer at current
time step



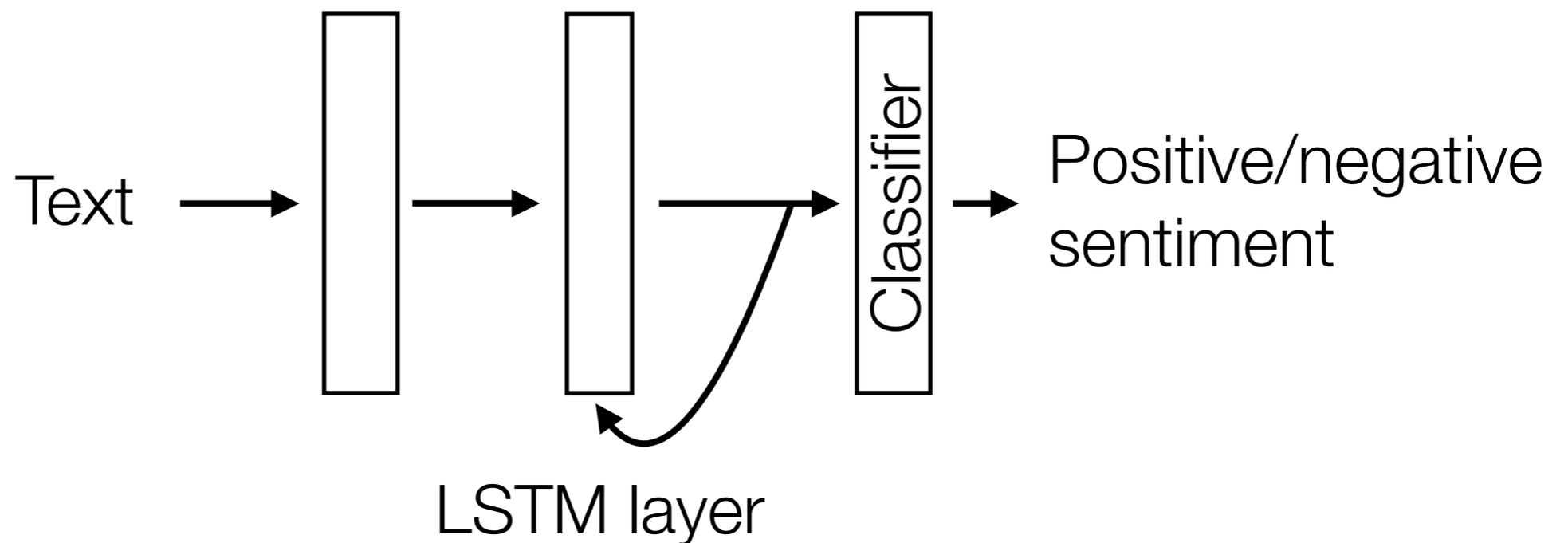
In `keras`, different
RNN options:
`SimpleRNN`, `LSTM`

like a dense layer
that has memory

Recommendation:
use LSTMs if you want
to have longer memory
(long range structure)

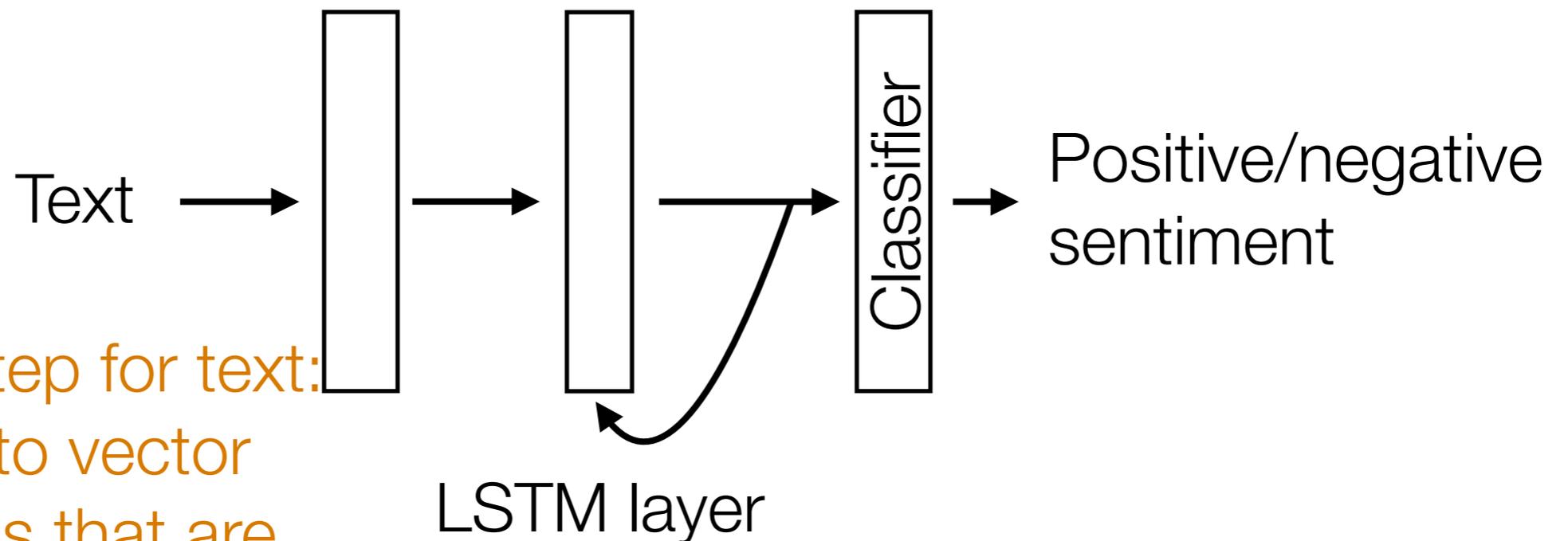
RNNs

Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)



RNNs

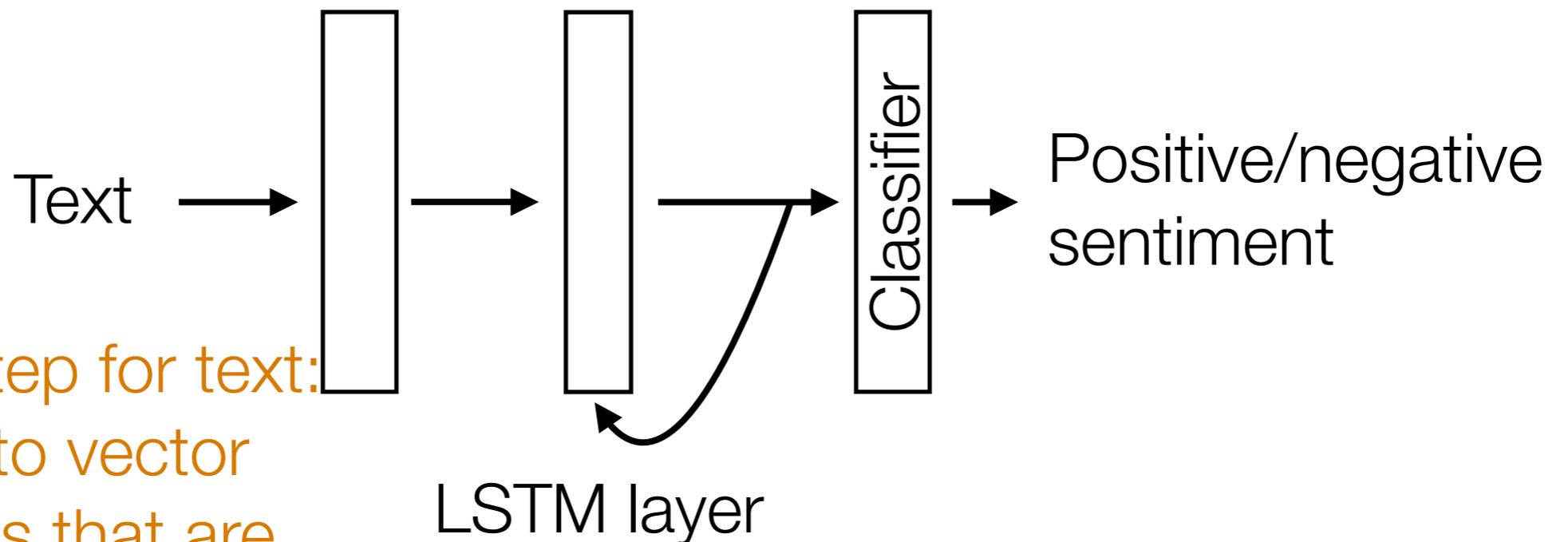
Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)



Common first step for text:
turn words into vector
representations that are
semantically meaningful

RNNs

Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)

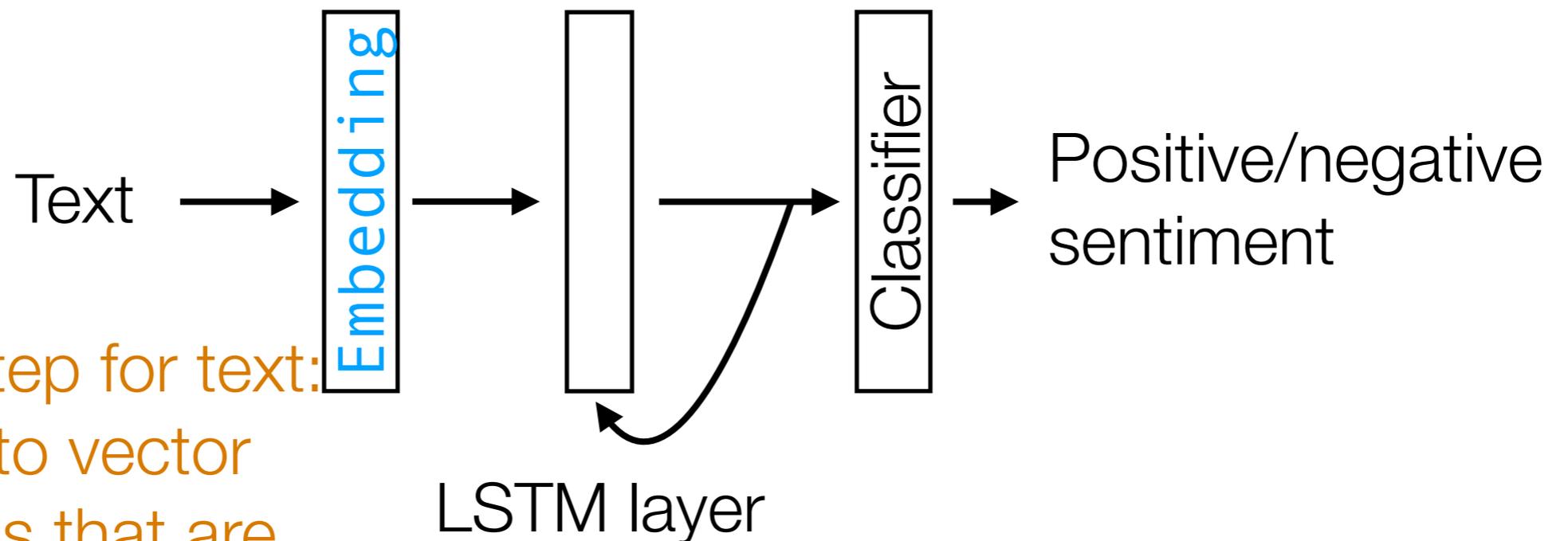


Common first step for text:
turn words into vector
representations that are
semantically meaningful

In `keras`, use
`Embedding` layer

RNNs

Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)

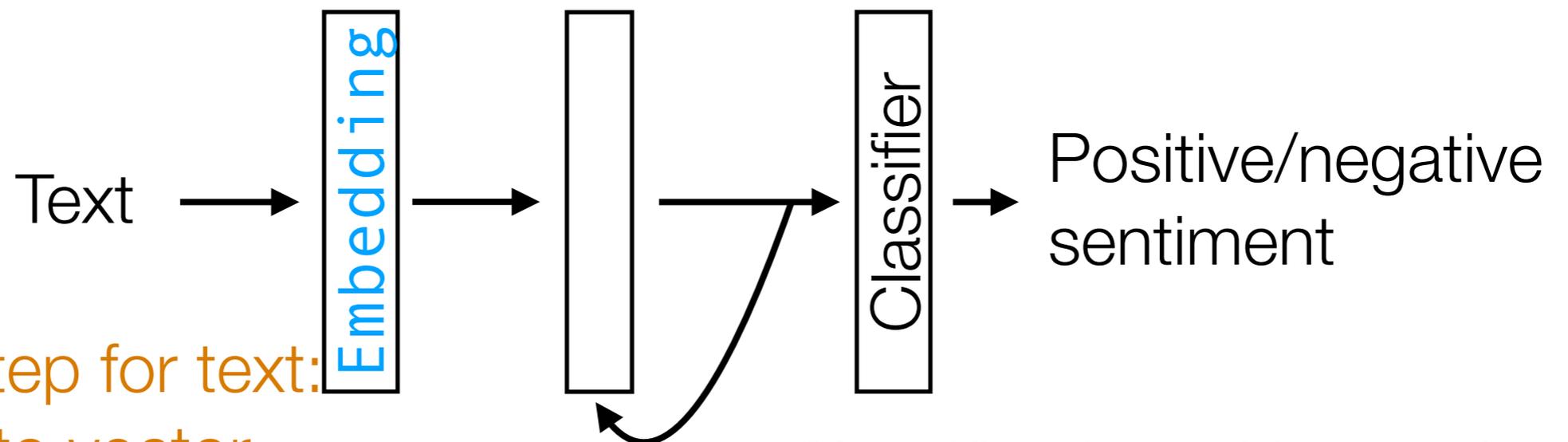


Common first step for text:
turn words into vector
representations that are
semantically meaningful

In `keras`, use
`Embedding` layer

RNNs

Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)



Common first step for text:
turn words into vector representations that are semantically meaningful

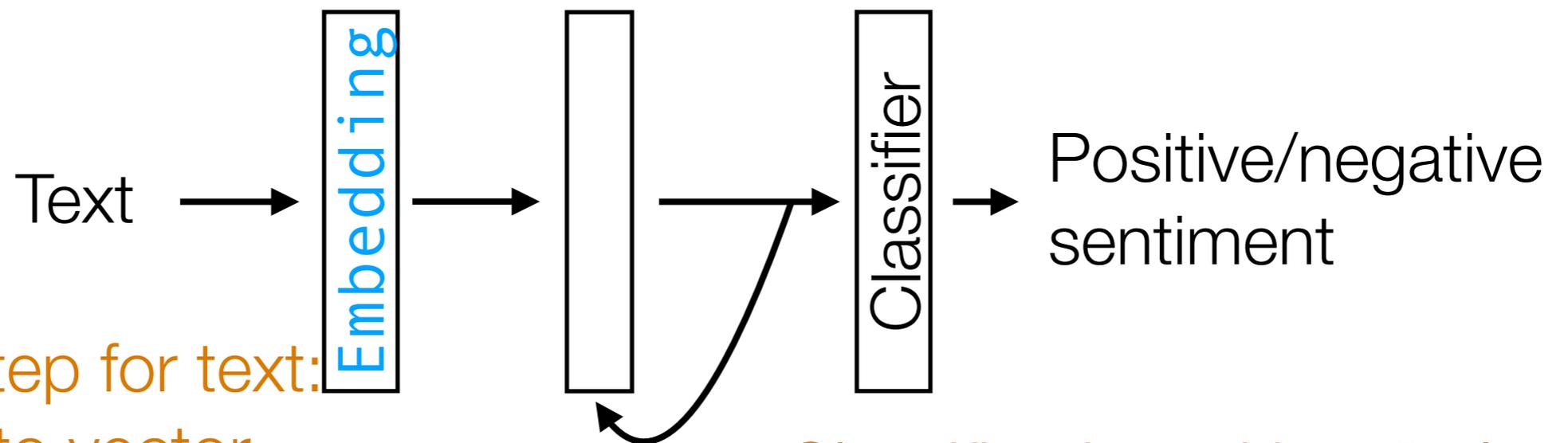
LSTM layer

Classification with > 2 classes:
dense layer, softmax activation

In `keras`, use
`Embedding` layer

RNNs

Example: Given text (e.g., movie review, Tweet), figure out whether it has positive or negative sentiment (binary classification)



Common first step for text:
turn words into vector representations that are semantically meaningful

In `keras`, use `Embedding` layer

LSTM layer

Classification with > 2 classes:
dense layer, softmax activation

Classification with 2 classes:
dense layer with 1 neuron,
sigmoid activation

RNNs

Demo

RNNs

RNNs

- Neatly handles time series in which there is some sort of global structure, so memory helps

RNNs

- Neatly handles time series in which there is some sort of global structure, so memory helps
- If time series doesn't actually have global structure, performance gain from using RNNs could be little compared to using 1D CNNs

RNNs

- Neatly handles time series in which there is some sort of global structure, so memory helps
 - If time series doesn't actually have global structure, performance gain from using RNNs could be little compared to using 1D CNNs
- An RNN layer should be chained together with other layers that learn a semantically meaningful interpretation from data (e.g., CNNs for images, word embeddings like word2vec/GloVe for text)