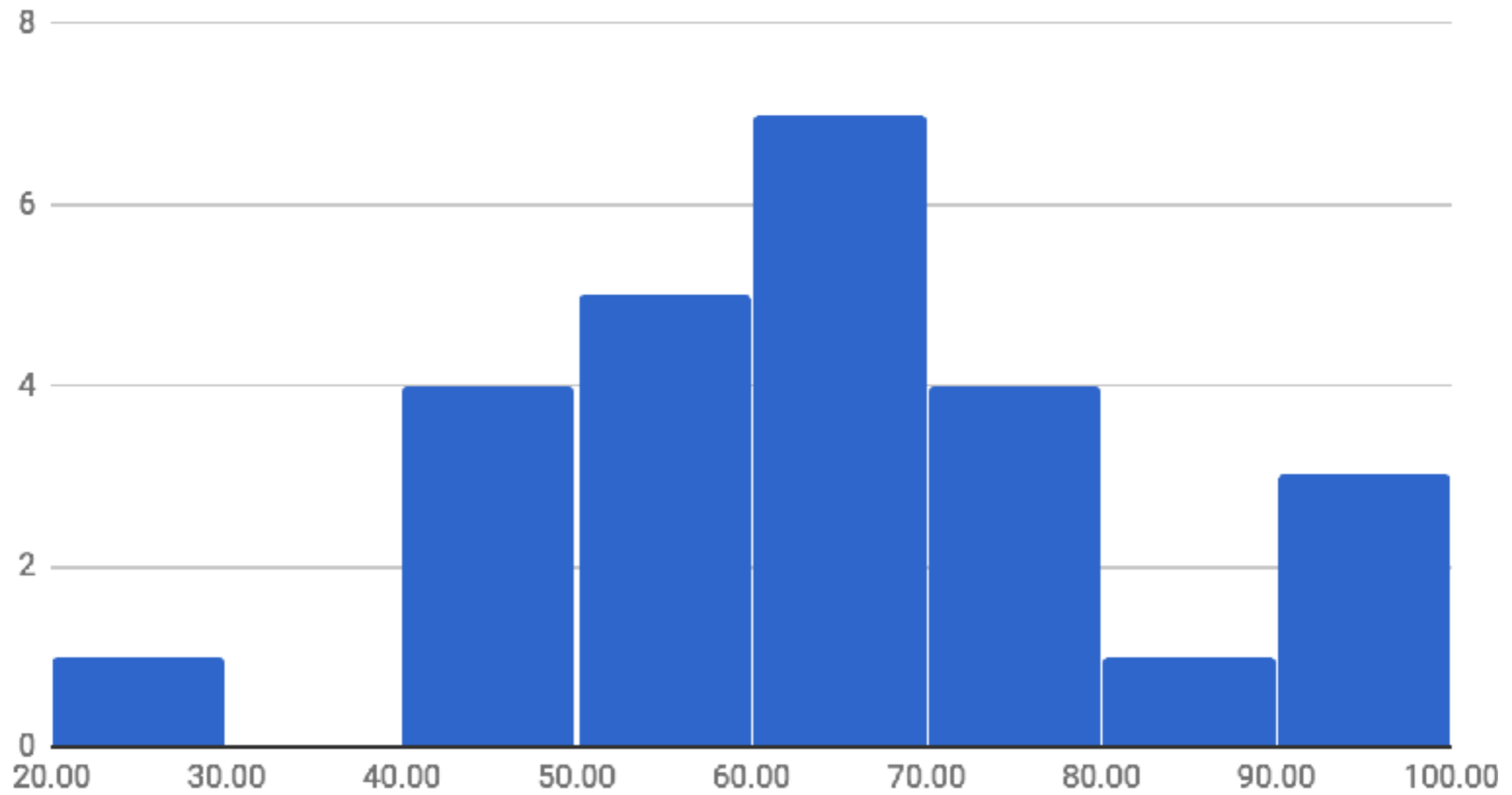


Image analysis with CNNs

George Chen

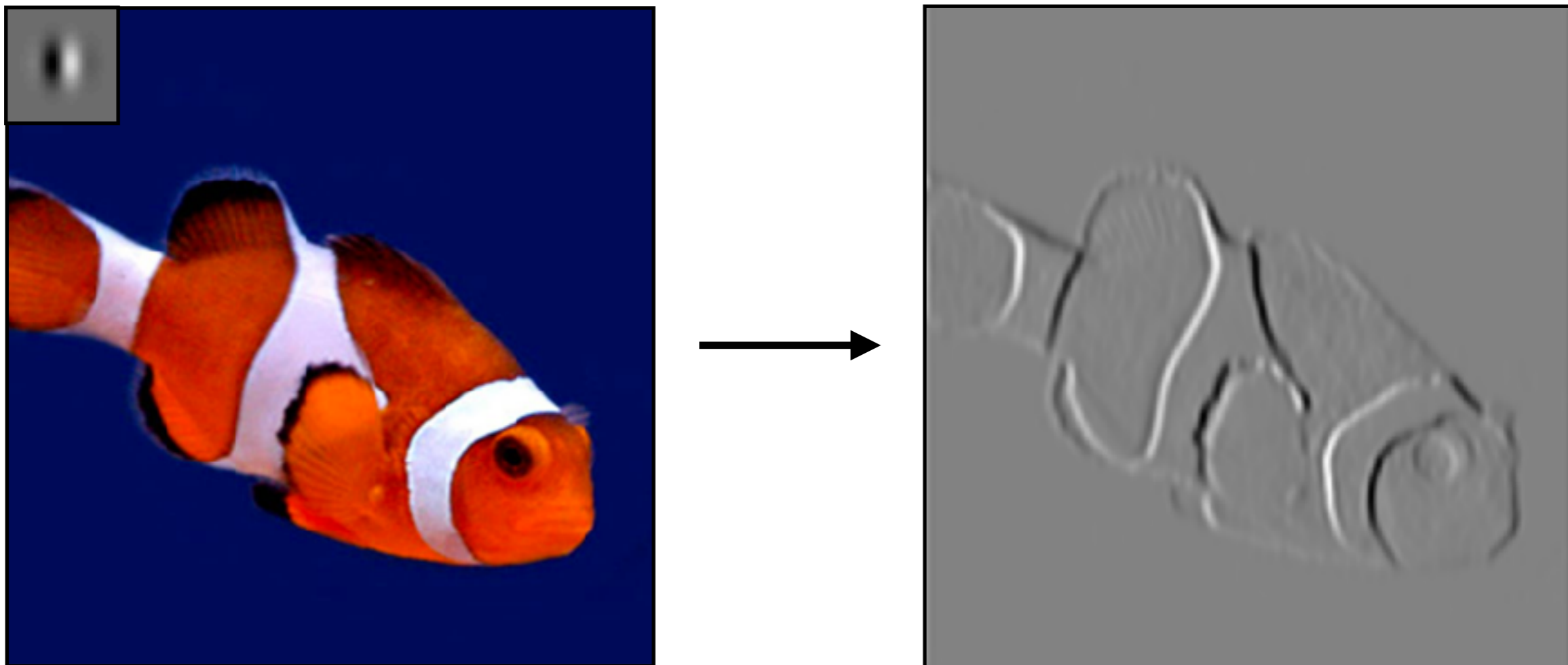
1 slide is by Phillip Isola

Mid-mini quiz histogram

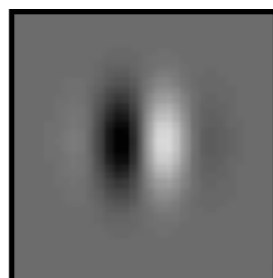


**Image analysis with
Convolutional Neural Nets
(CNNs, also called convnets)**

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

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!

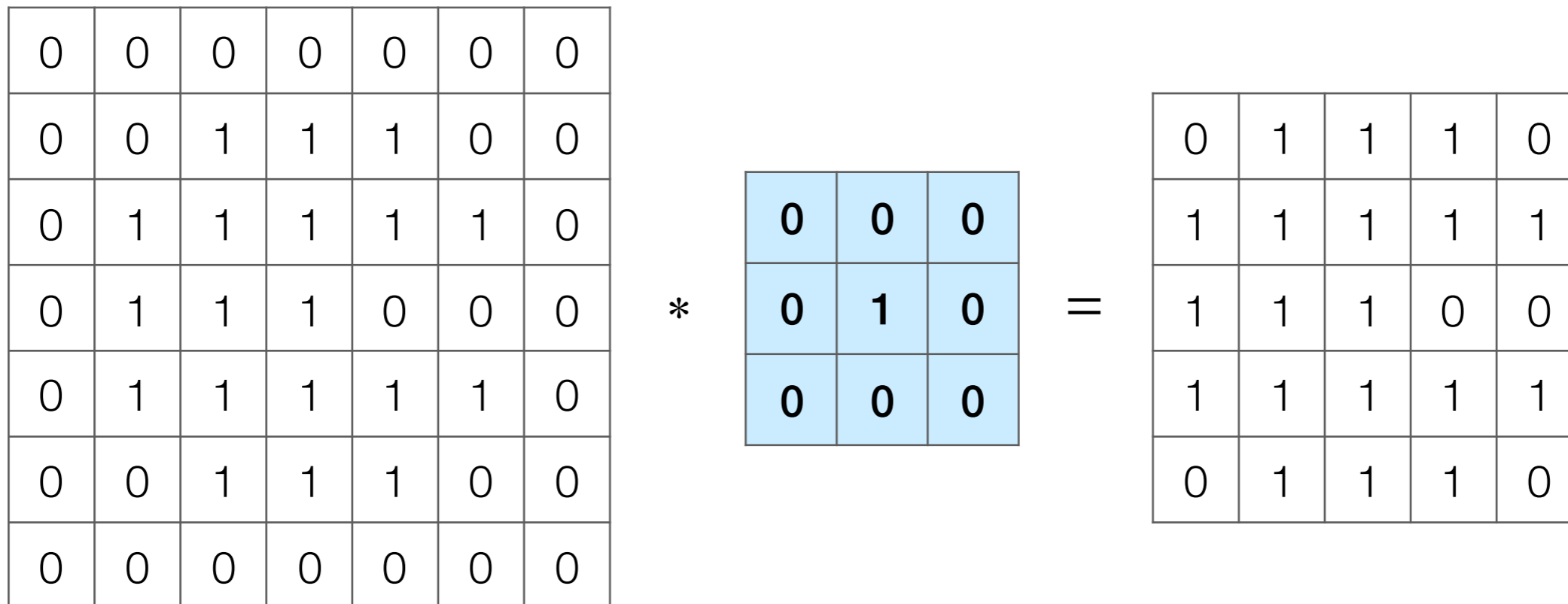
0	0	0	0	0	0	0	
0	0	1	1	0	1	0	0
0	1	1	1	0	1	1	0
0	1	1	0	1	0	0	0
0	1	1	1	1	1	1	0
0	0	1	1	1	0	0	
0	0	0	0	0	0	0	

Input image

0	1	1	1	0
1	1			

Output image

Convolution



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	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

*	$\frac{1}{9}$	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

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

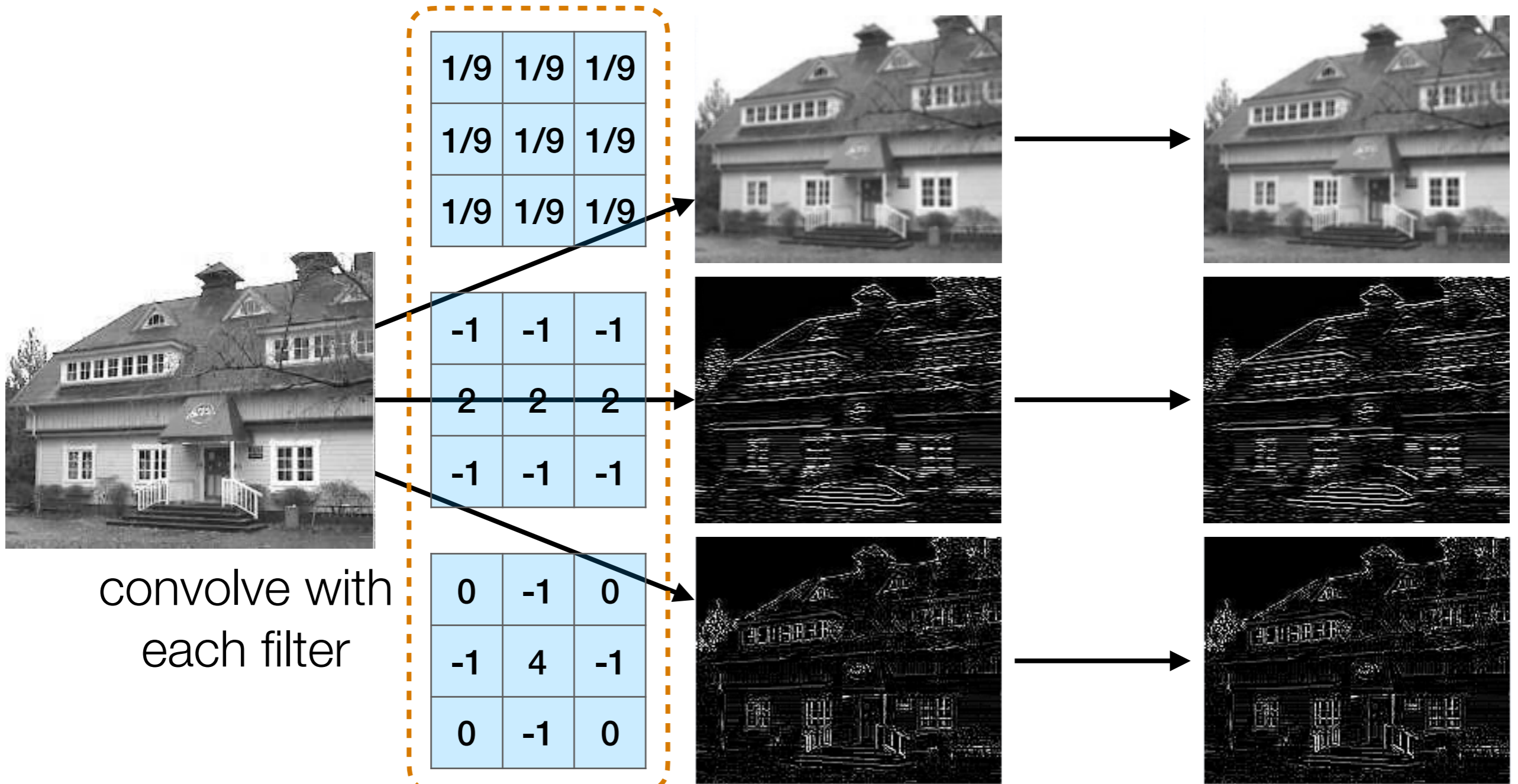


$$\begin{matrix} * & \begin{matrix} \begin{matrix} -1 & -1 & -1 \\ 2 & 2 & 2 \\ -1 & -1 & -1 \end{matrix} \\ = \end{matrix} \end{matrix}$$



(this example finds horizontal edges)

Convolution Layer

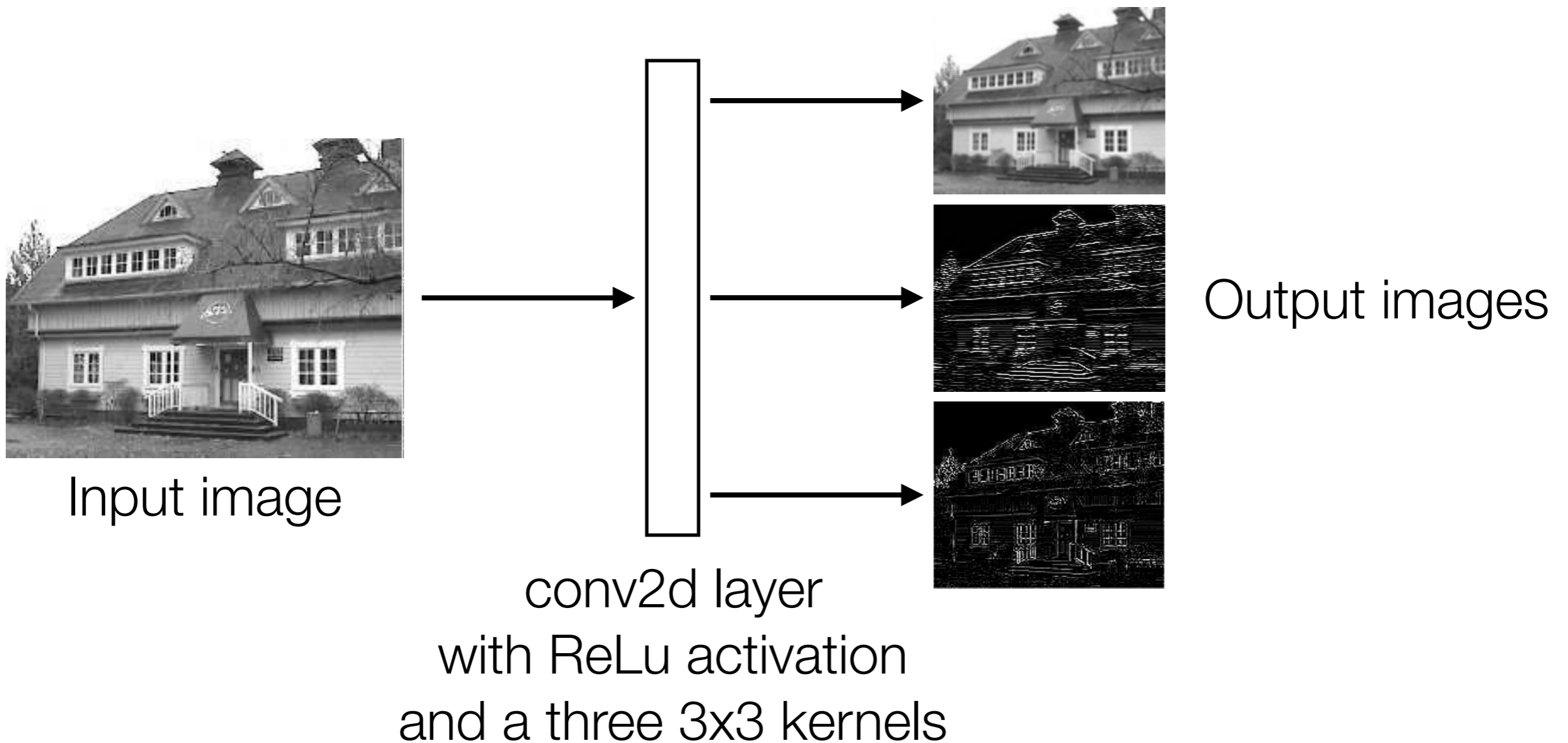


convolve with each filter

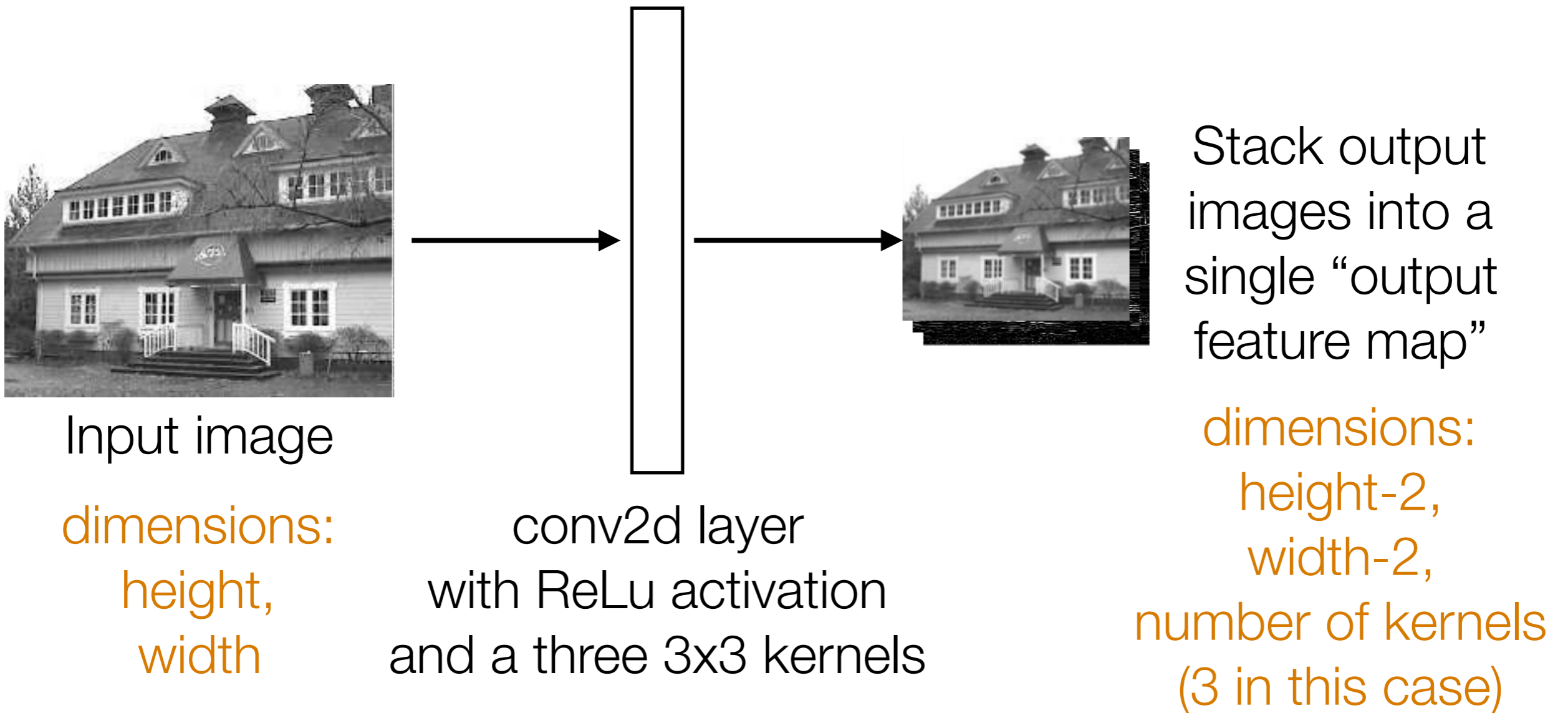
filters are actually unknown and are learned!

activation (e.g., ReLU)

Convolution Layer



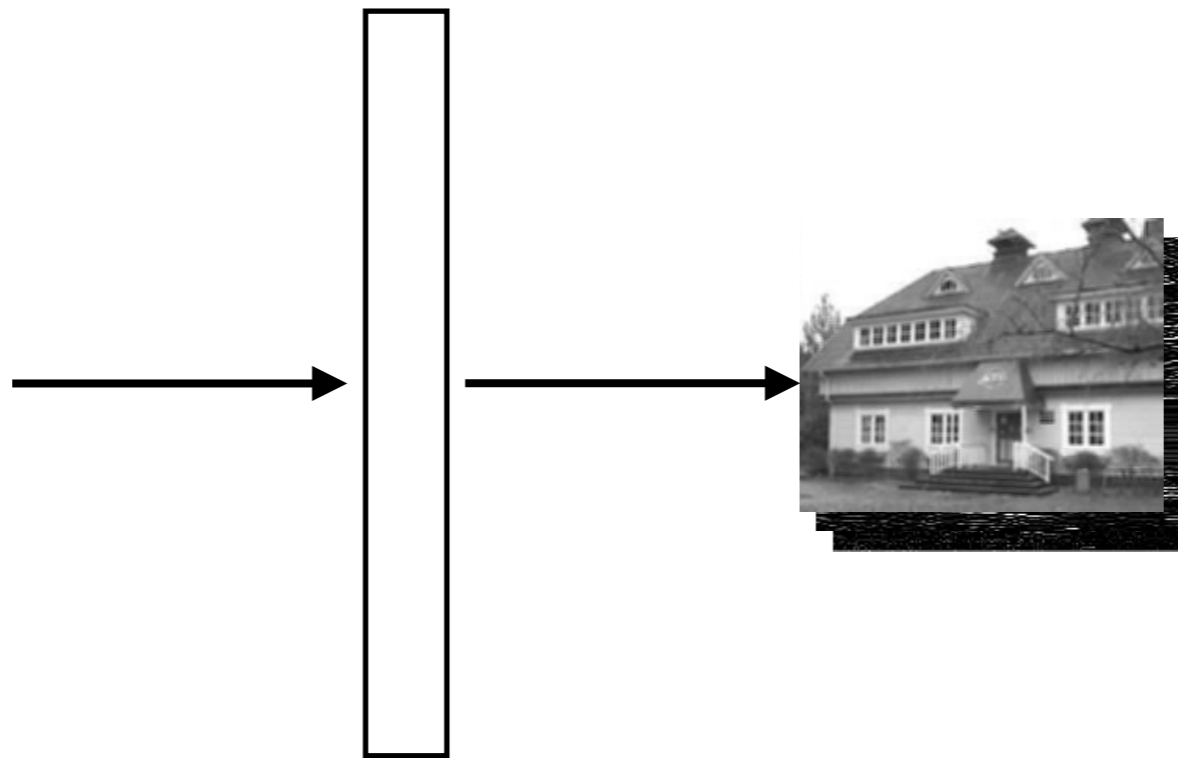
Convolution Layer



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

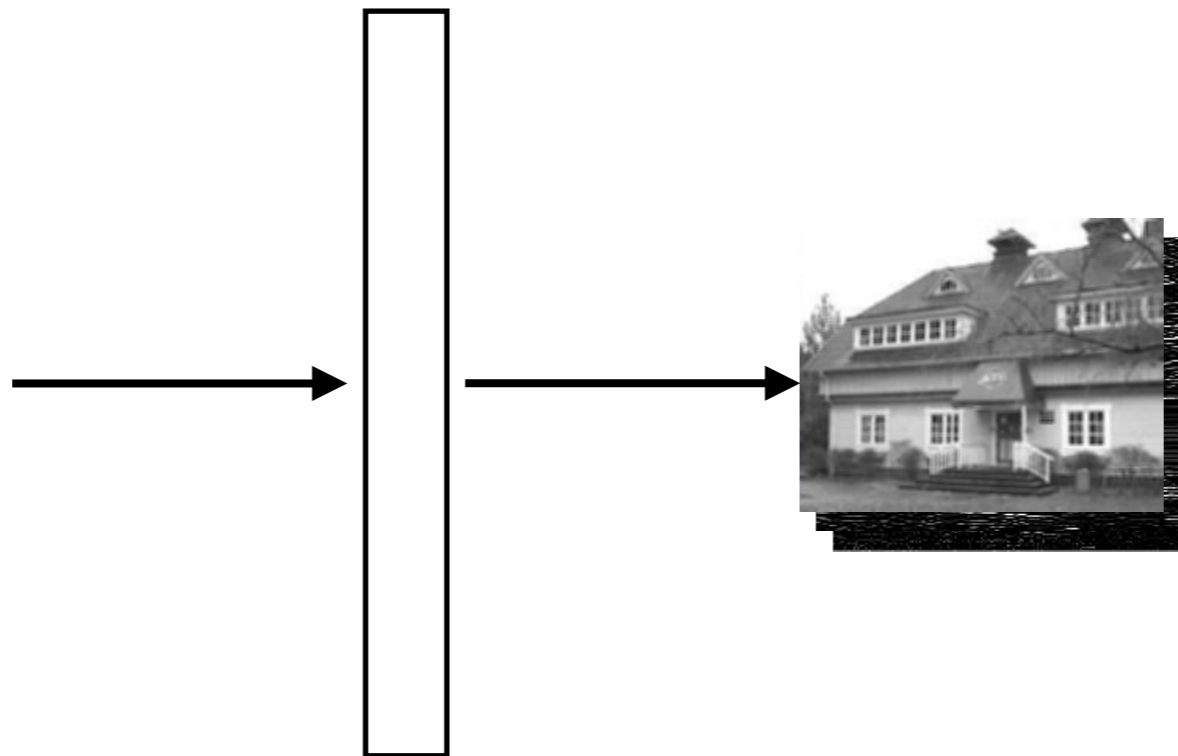
Convolution Layer



Input image

dimensions:
height,
width,

depth d (# channels)



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

technical detail: there's
also a bias vector



Stack output
images into a
single “output
feature map”

dimensions:
height-2,
width-2,
 k

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

	-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

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

What numbers were involved in computing this 1?

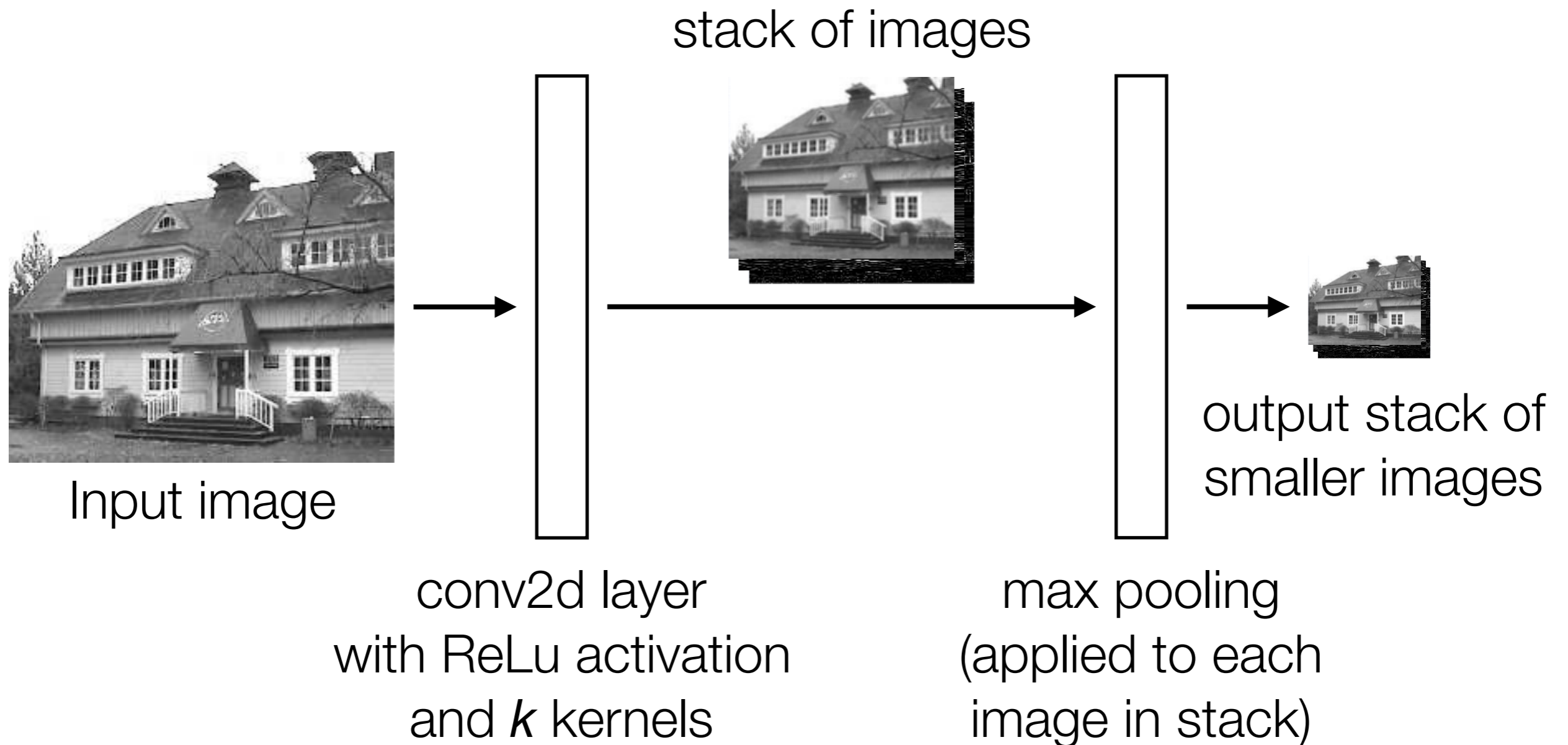
In this example: 1 pixel in max pooling output captures information from 16 input pixels!

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

1	3
1	3

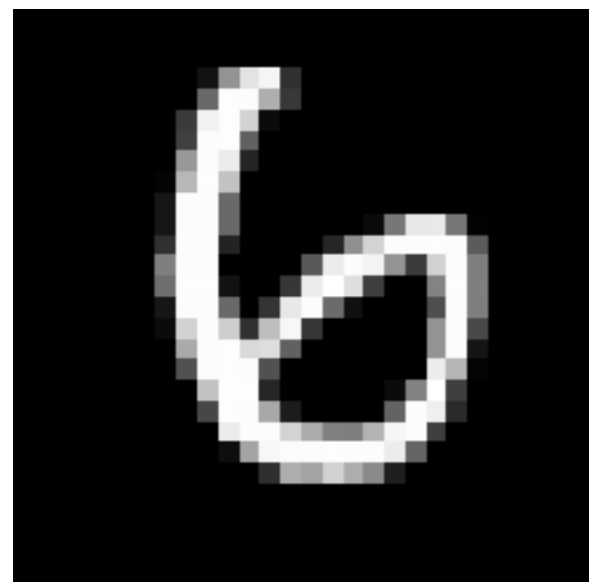
Output after max pooling

Basic Building Block of CNN's



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**

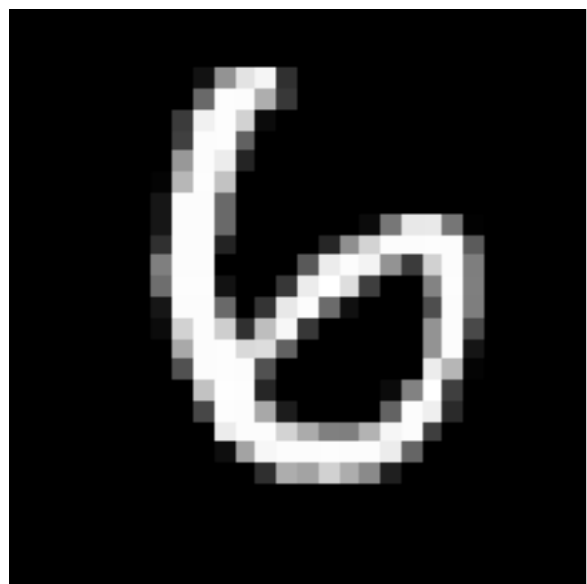
$$\log \frac{1}{\text{Pr}(\text{digit } 6)}$$

Error is averaged across training examples

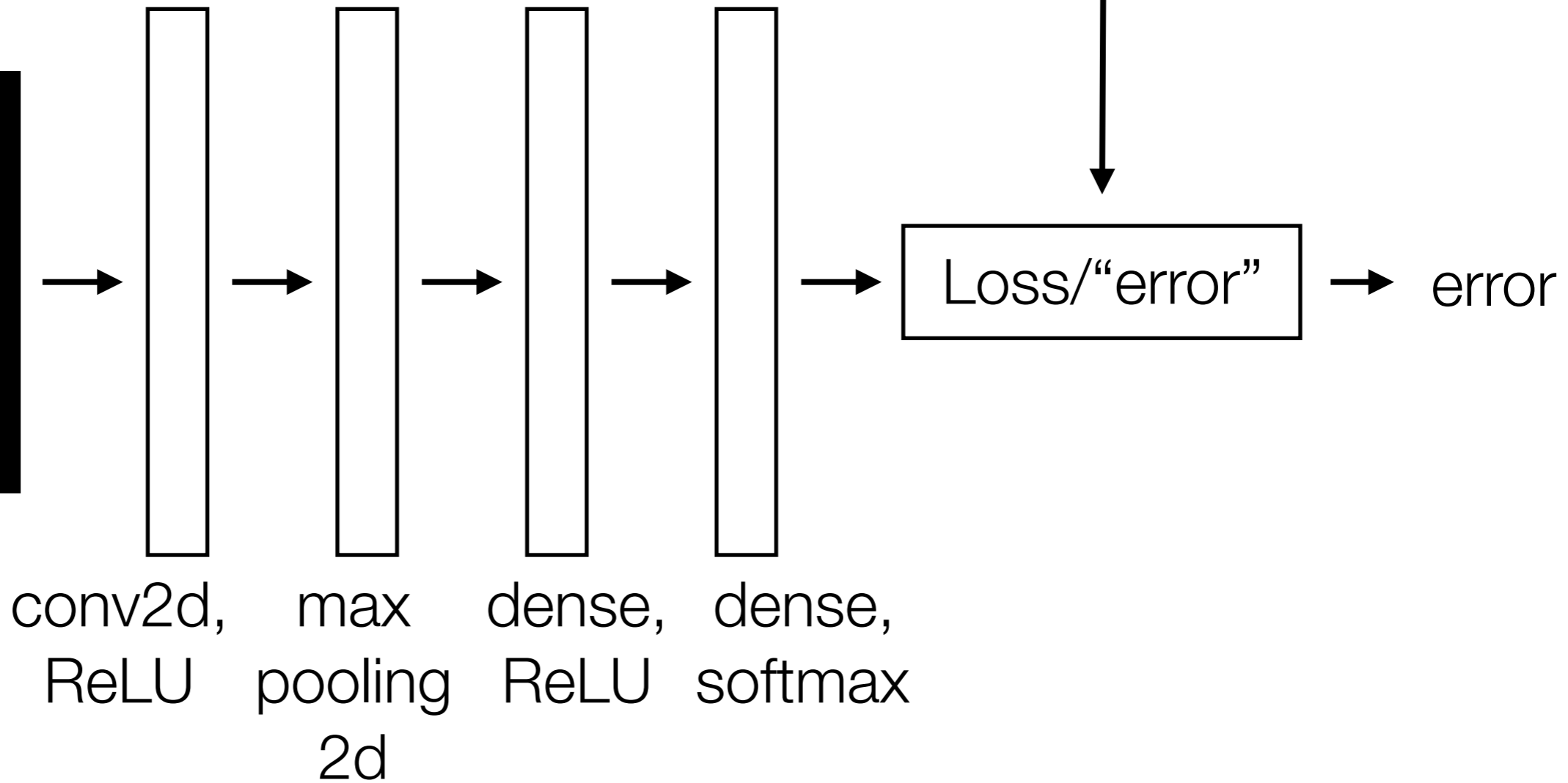
error

Handwritten Digit Recognition

Training label: 6



28x28 image

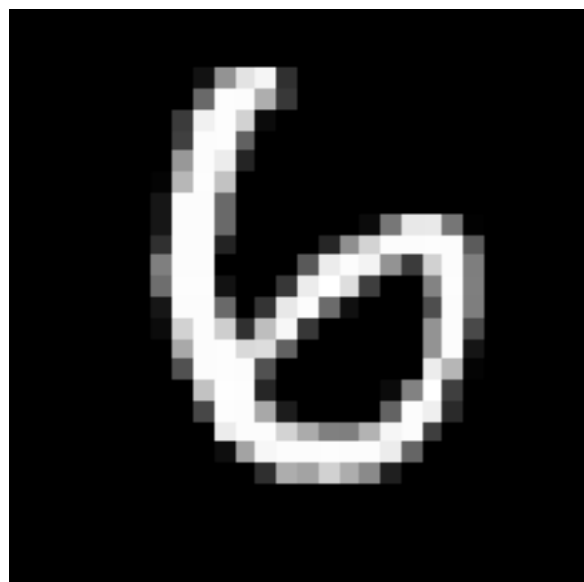


Handwritten Digit Recognition

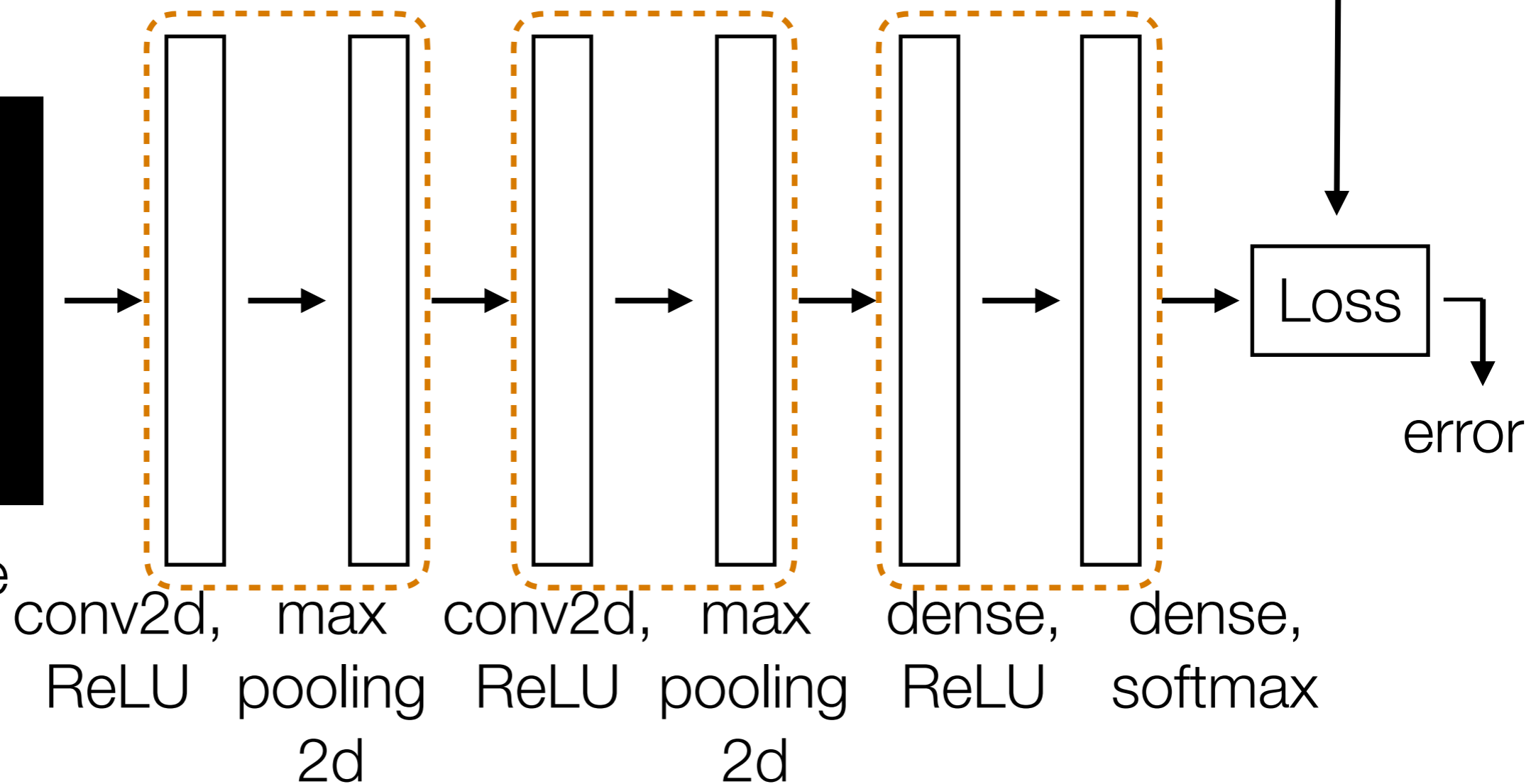
Training label: 6

extract low-level visual features & aggregate

non-vision-specific classification neural net



28x28 image



extract higher-level visual features & aggregate

CNN Demo

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