

SSY098 - Image Analysis

Lecture 6 - (More) Convolutional Neural Networks

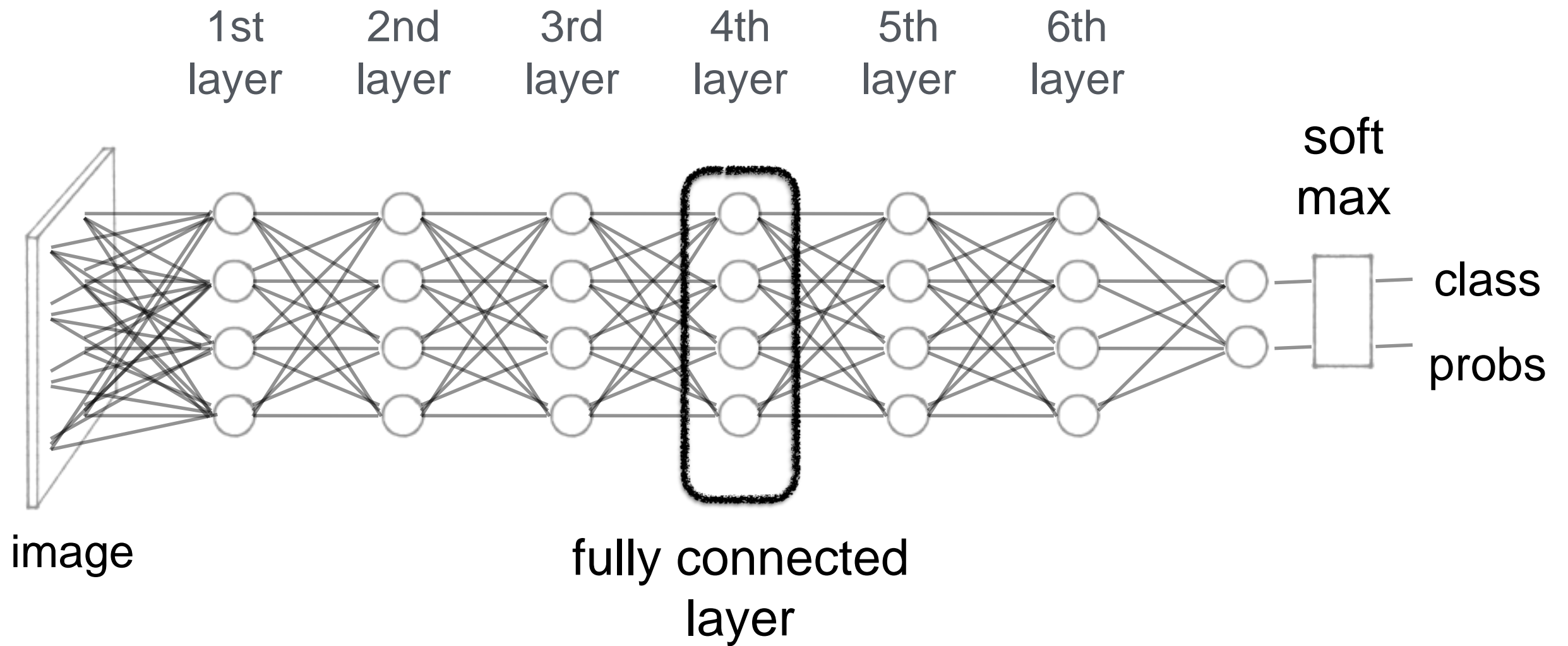
Torsten Sattler

(slides adapted from Olof Enqvist)

Last Lecture

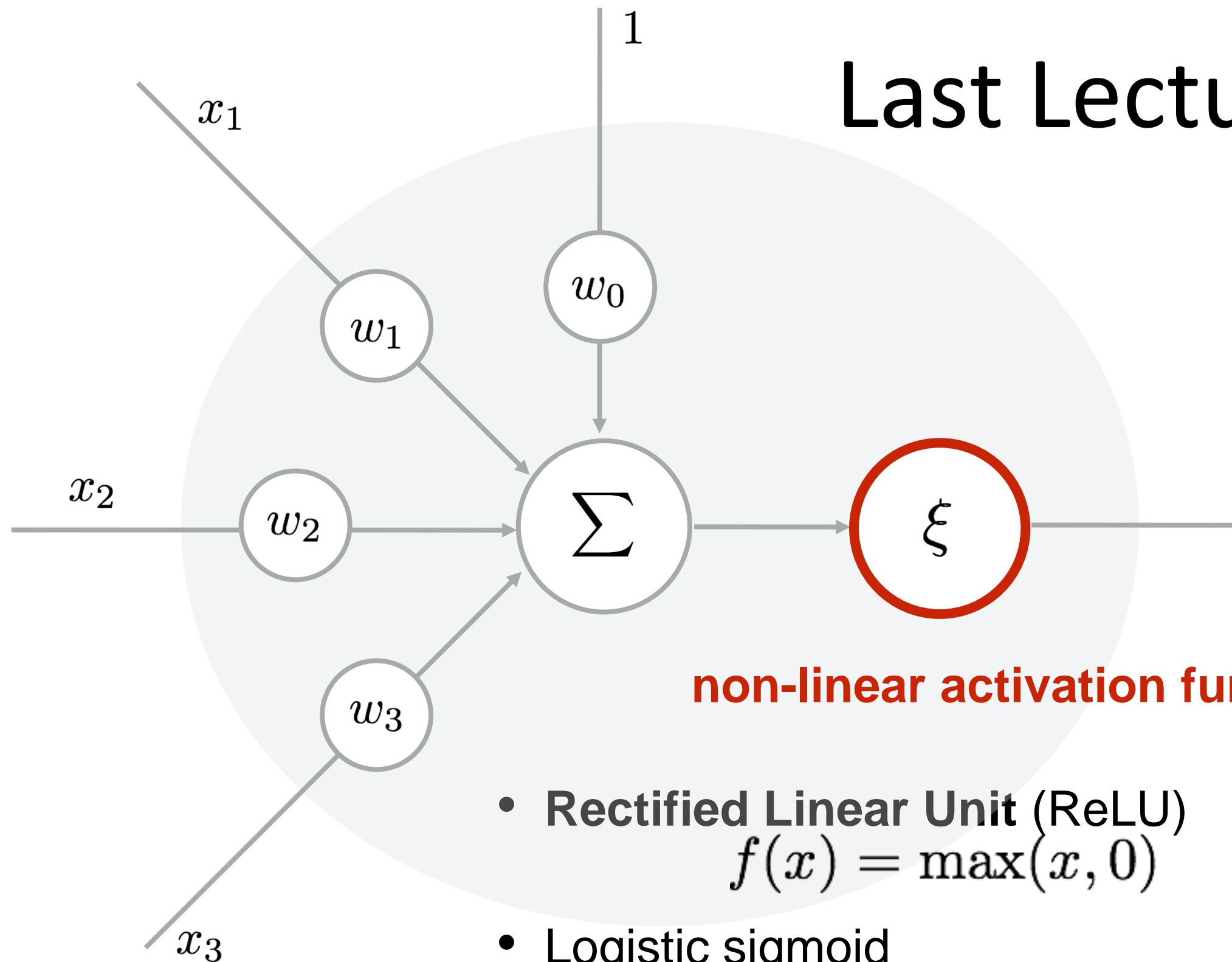
Jan. 20	Introduction, Linear classifiers and filtering	
Jan. 23	Filtering, gradients, scale	Lab 1
Jan. 27	Local features	
Jan. 30	Learning a classifier	Lab 2
Feb. 3	Convolutional neural networks	
Feb. 6	More convolutional neural networks	
Feb. 10	Robust model fitting and RANSAC	Lab 3
Feb. 13	Image registration	
Feb. 17	Camera Geometry	Lab 4
Feb. 20	More camera geometry	
Feb. 24	Generative neural networks	
Feb. 27	Generative neural networks	
Mar. 2	TBA	
Mar. 9	TBA	

Last Lecture



Neural Networks

Last Lecture



non-linear activation function

- **Rectified Linear Unit (ReLU)**
 $f(x) = \max(x, 0)$
- Logistic sigmoid
- Hyperbolic tangent (tanh)
- ...

A neuron

Last Lecture

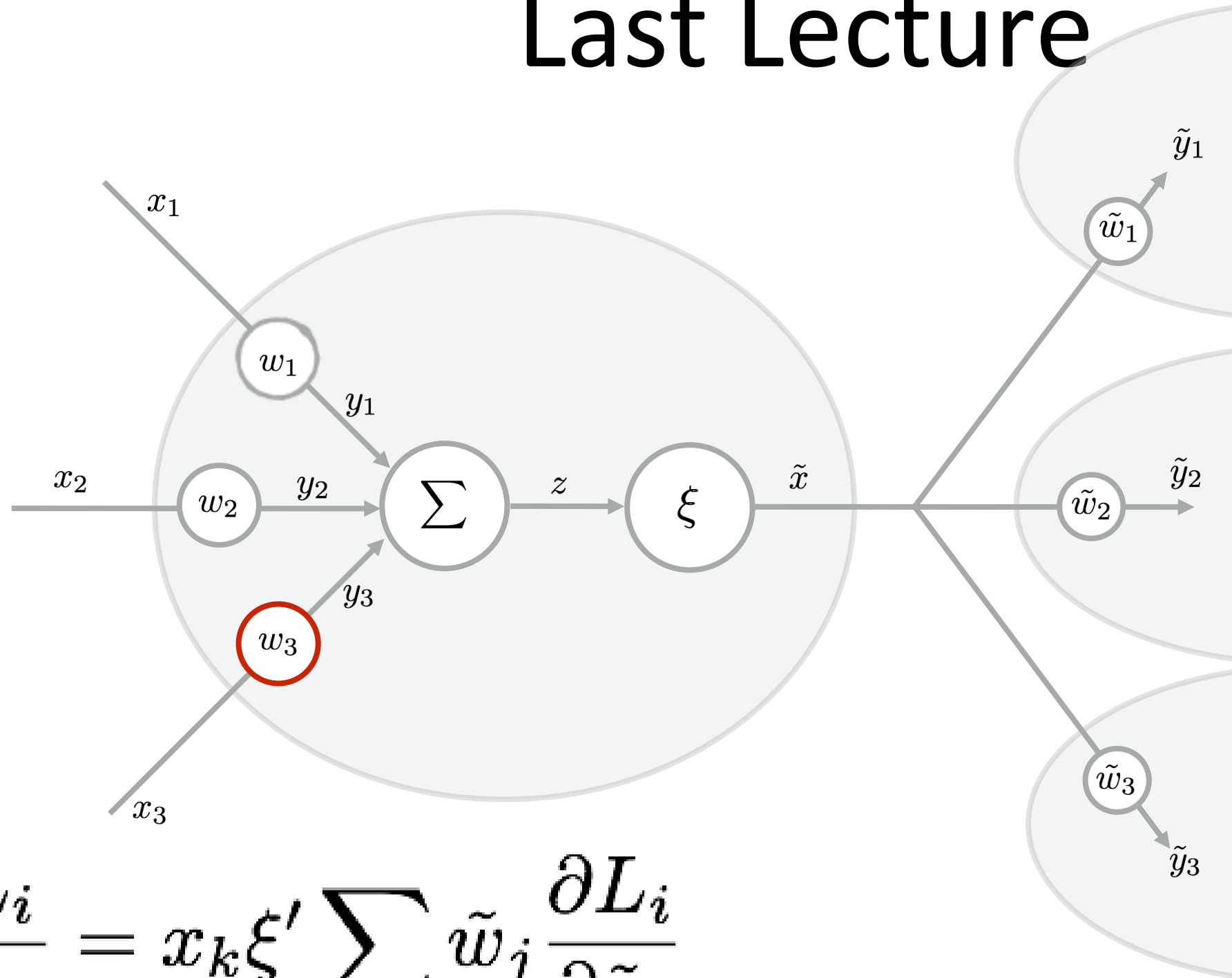
- Optimize all network parameters via gradient descent:

$$\theta^{(k+1)} = \theta^{(k)} - \mu \sum_i \nabla L_i(\theta) \approx \theta^{(k)} - \mu \nabla L_i(\theta)$$

- Compute derivative of partial loss for each parameter w_k

Training a neural network

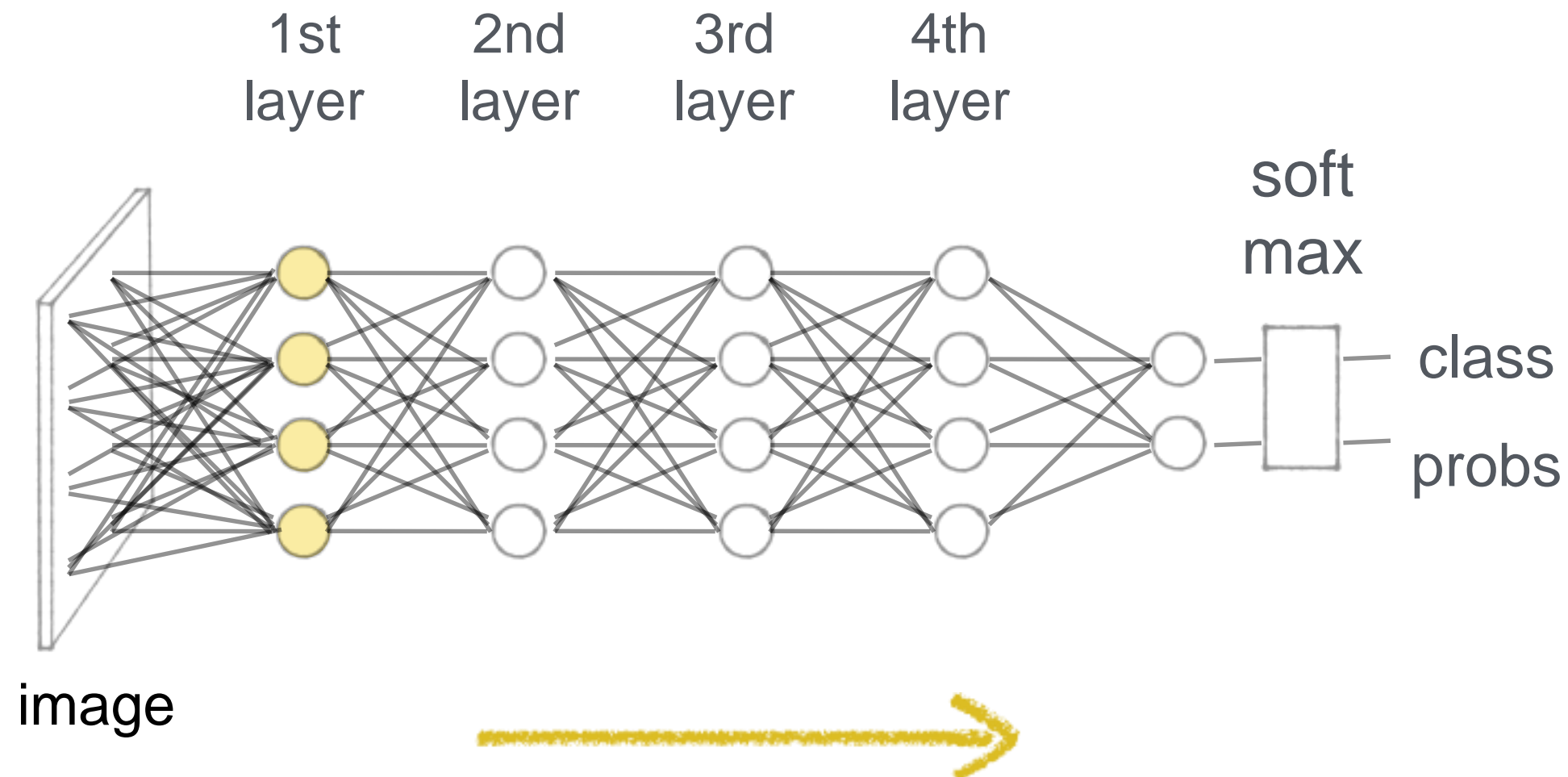
Last Lecture



$$\frac{\partial L_i}{\partial w_k} = x_k \xi' \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

Backpropagation (backprop)

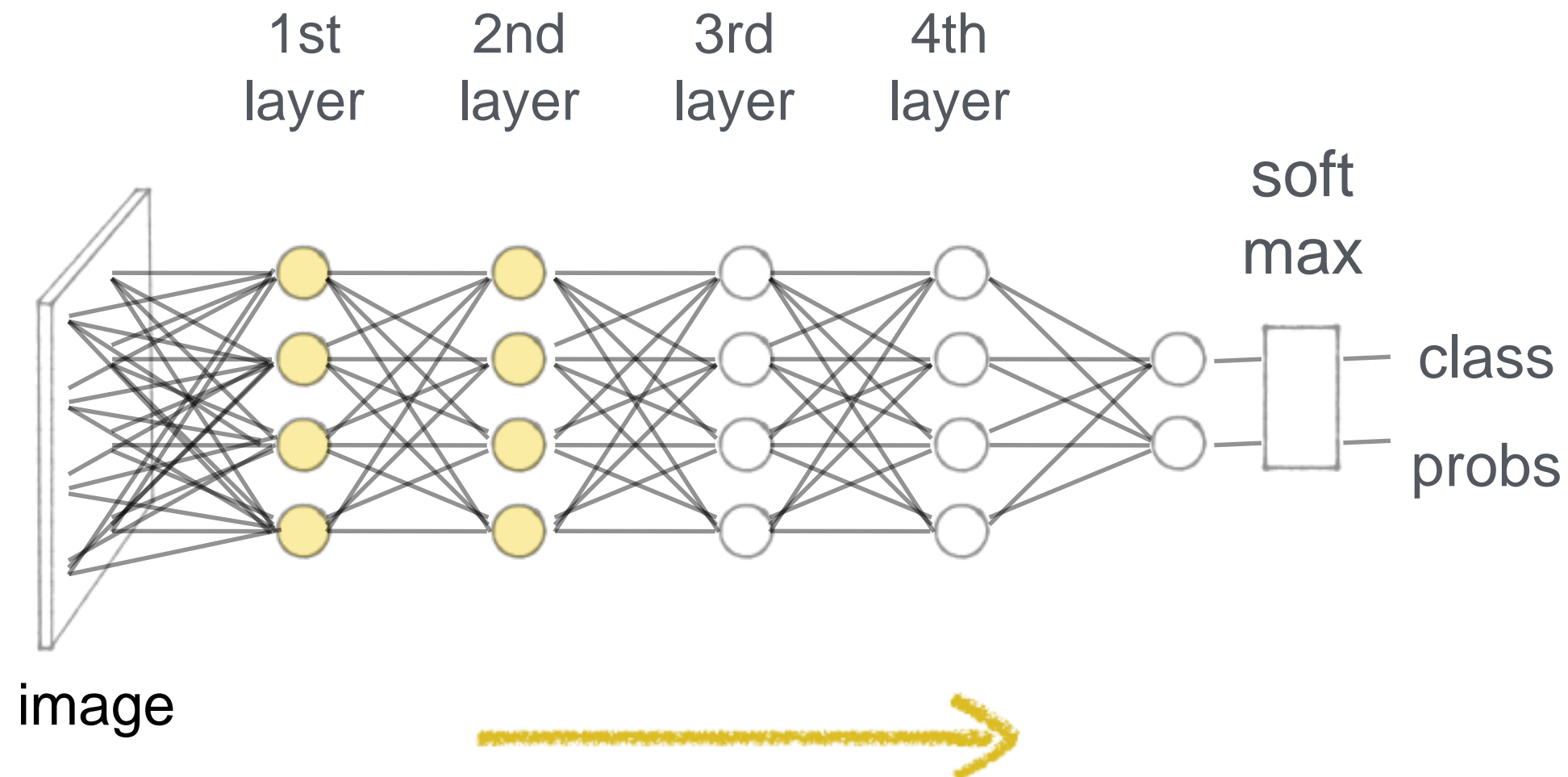
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

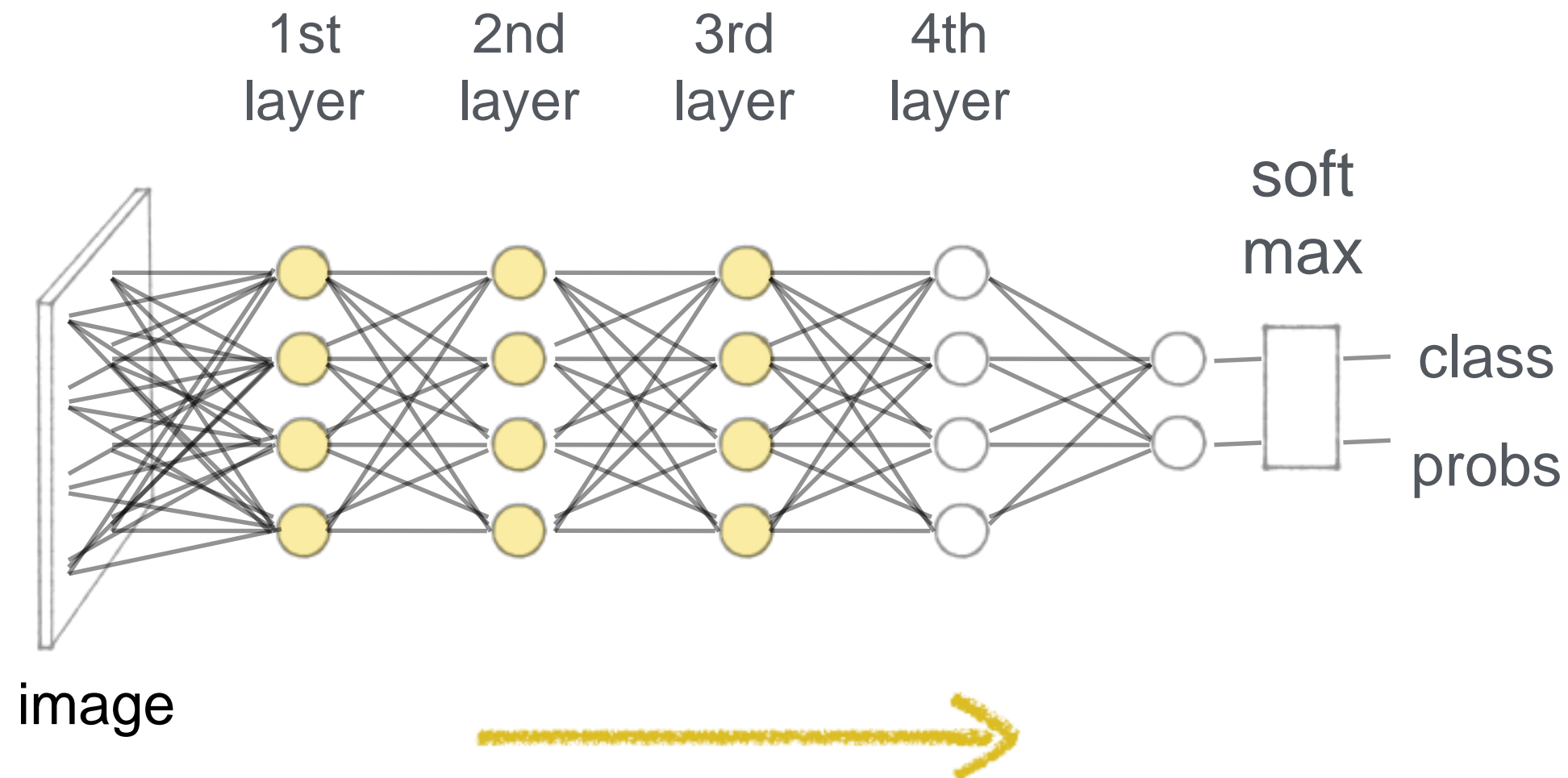
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

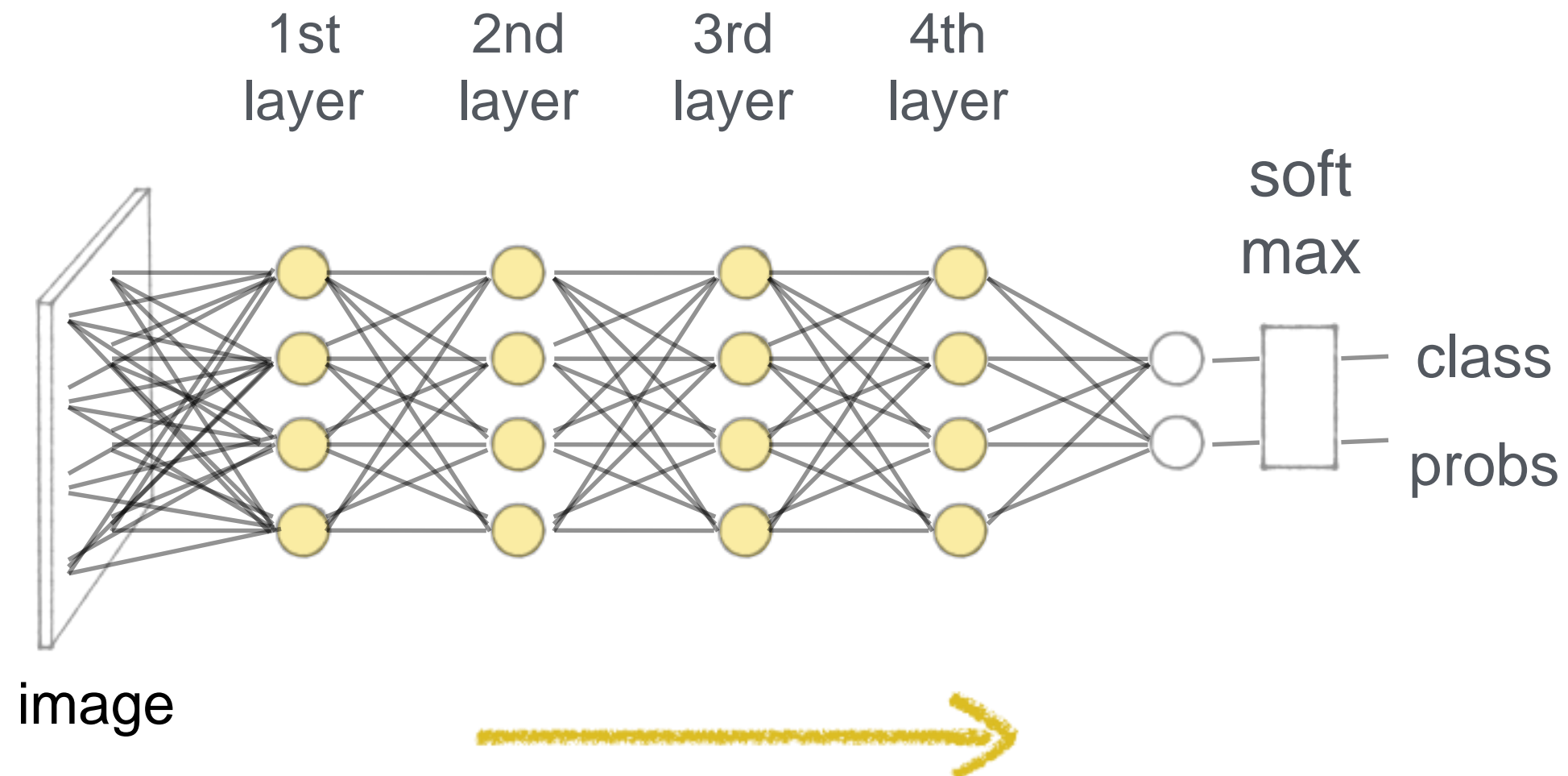
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

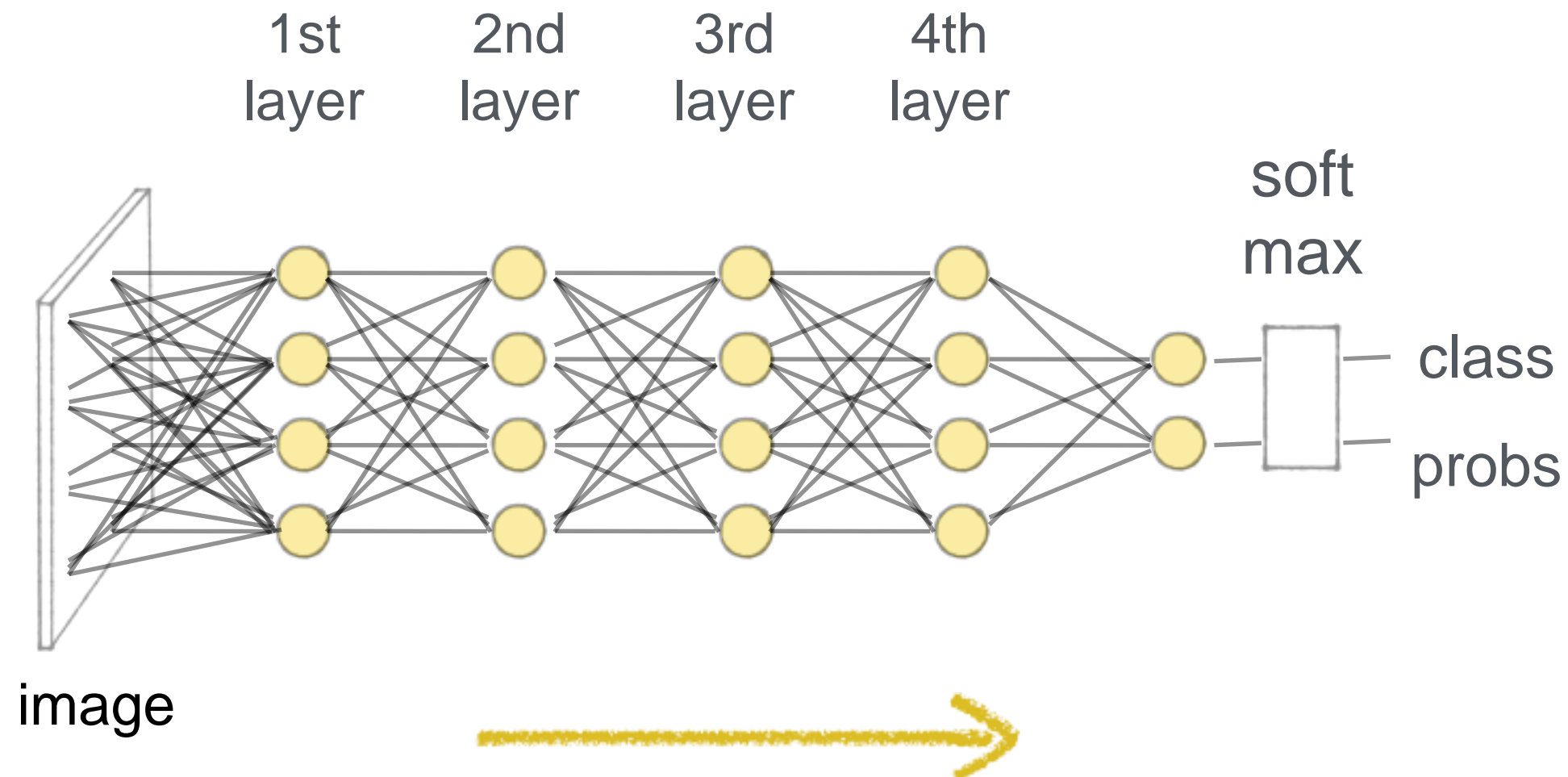
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

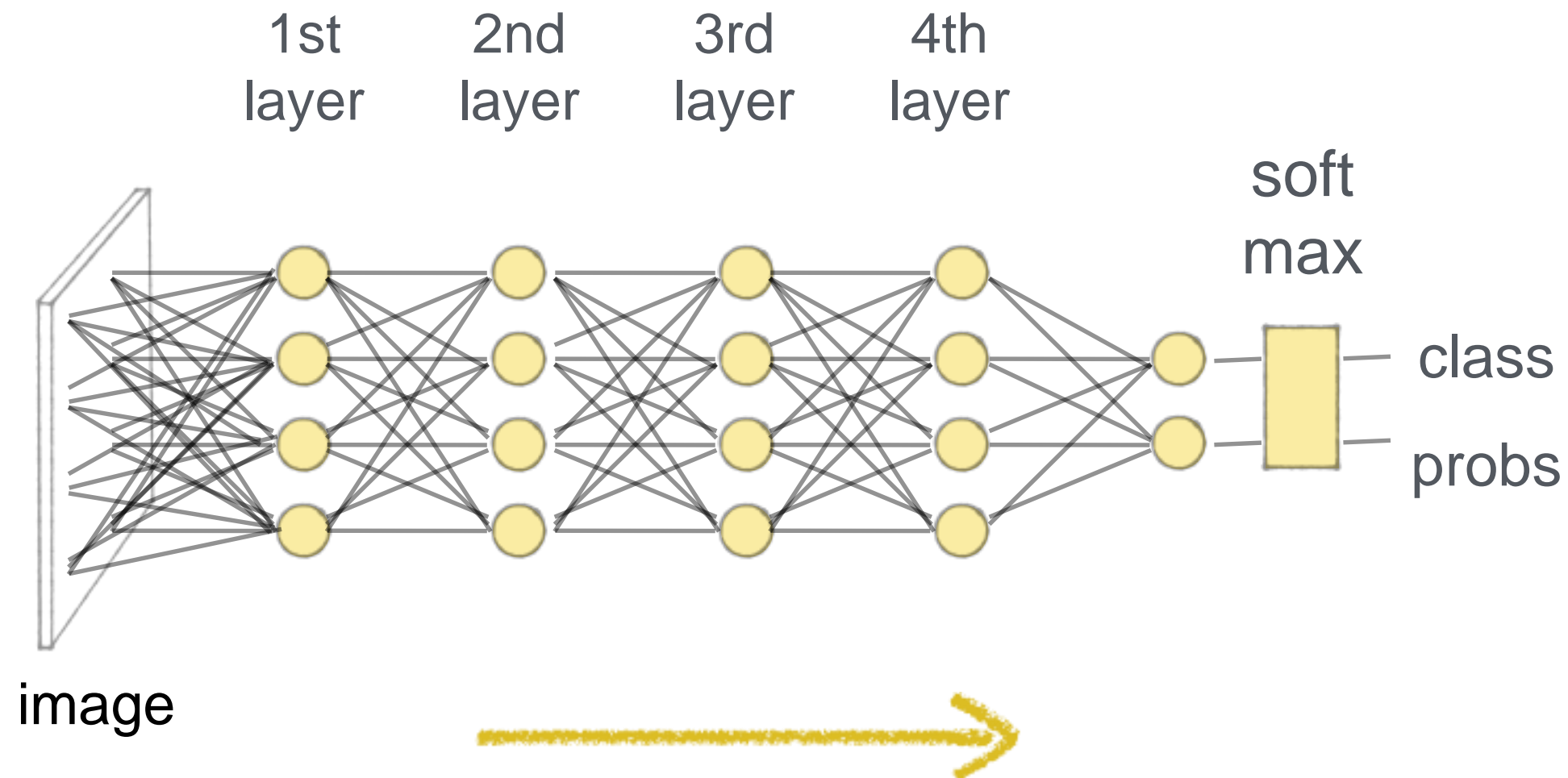
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

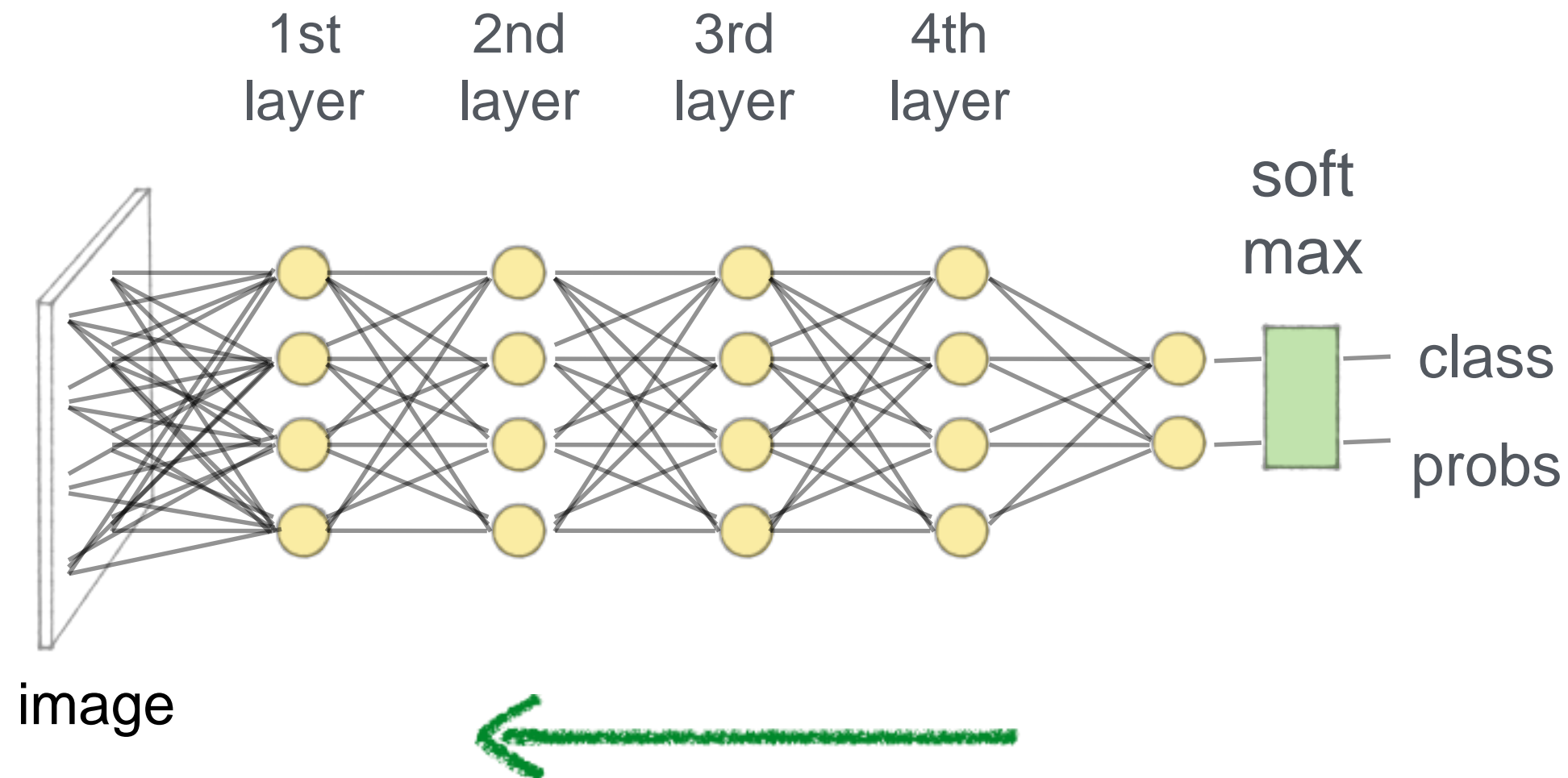
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underline{x_k \xi'} \sum_j \tilde{w}_j \frac{\partial L_i}{\partial \tilde{y}_j}$$

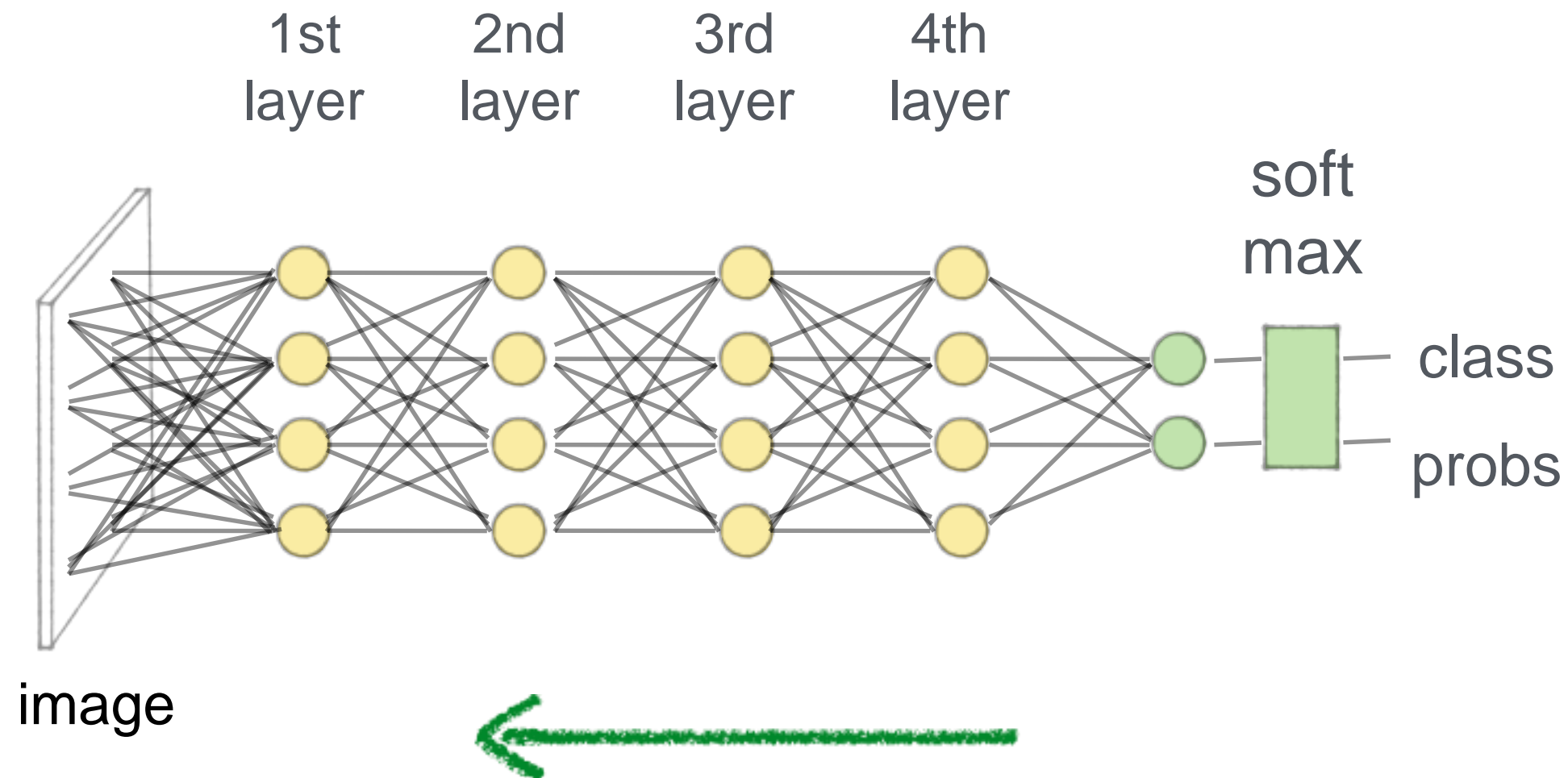
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

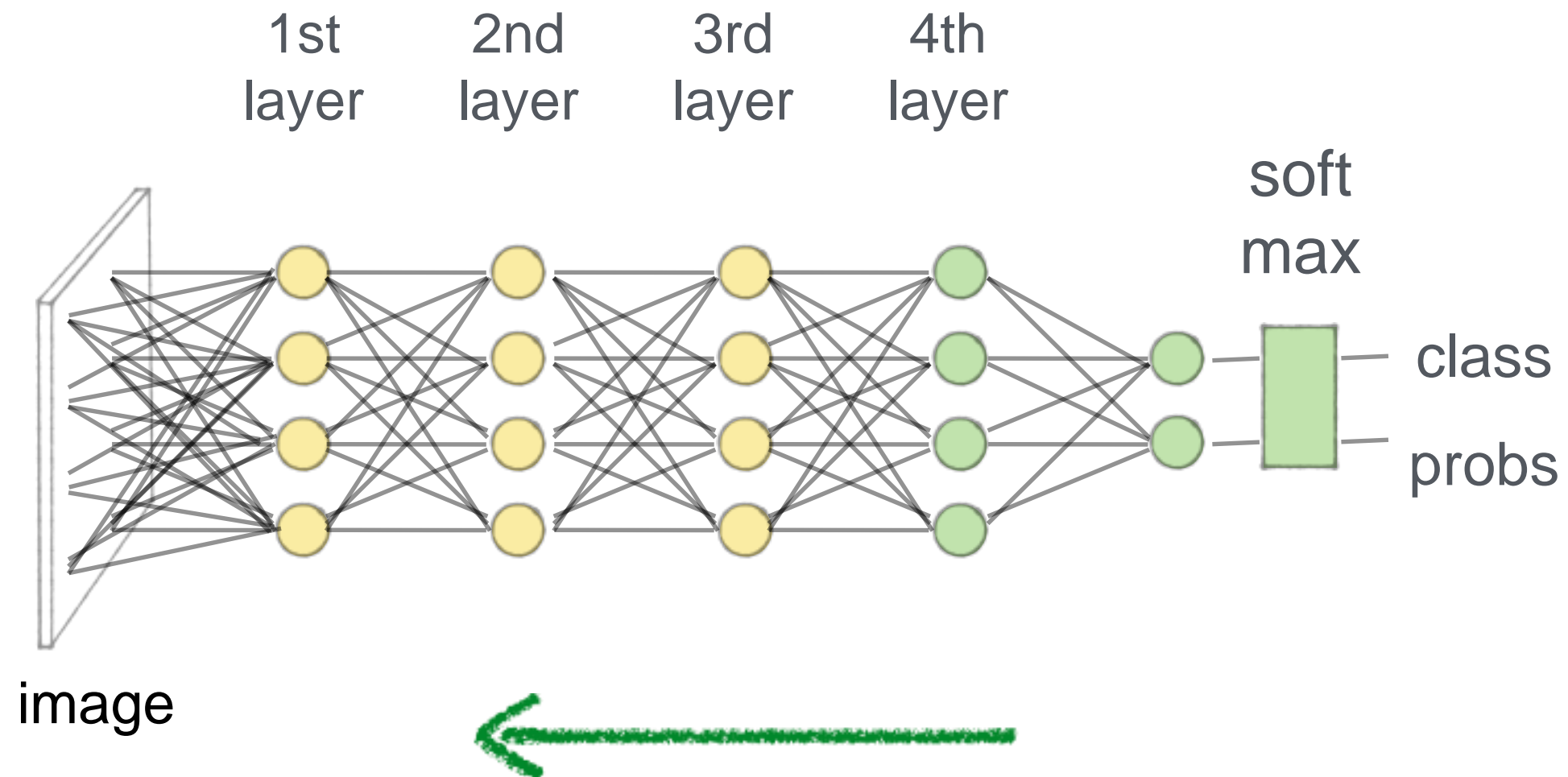
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

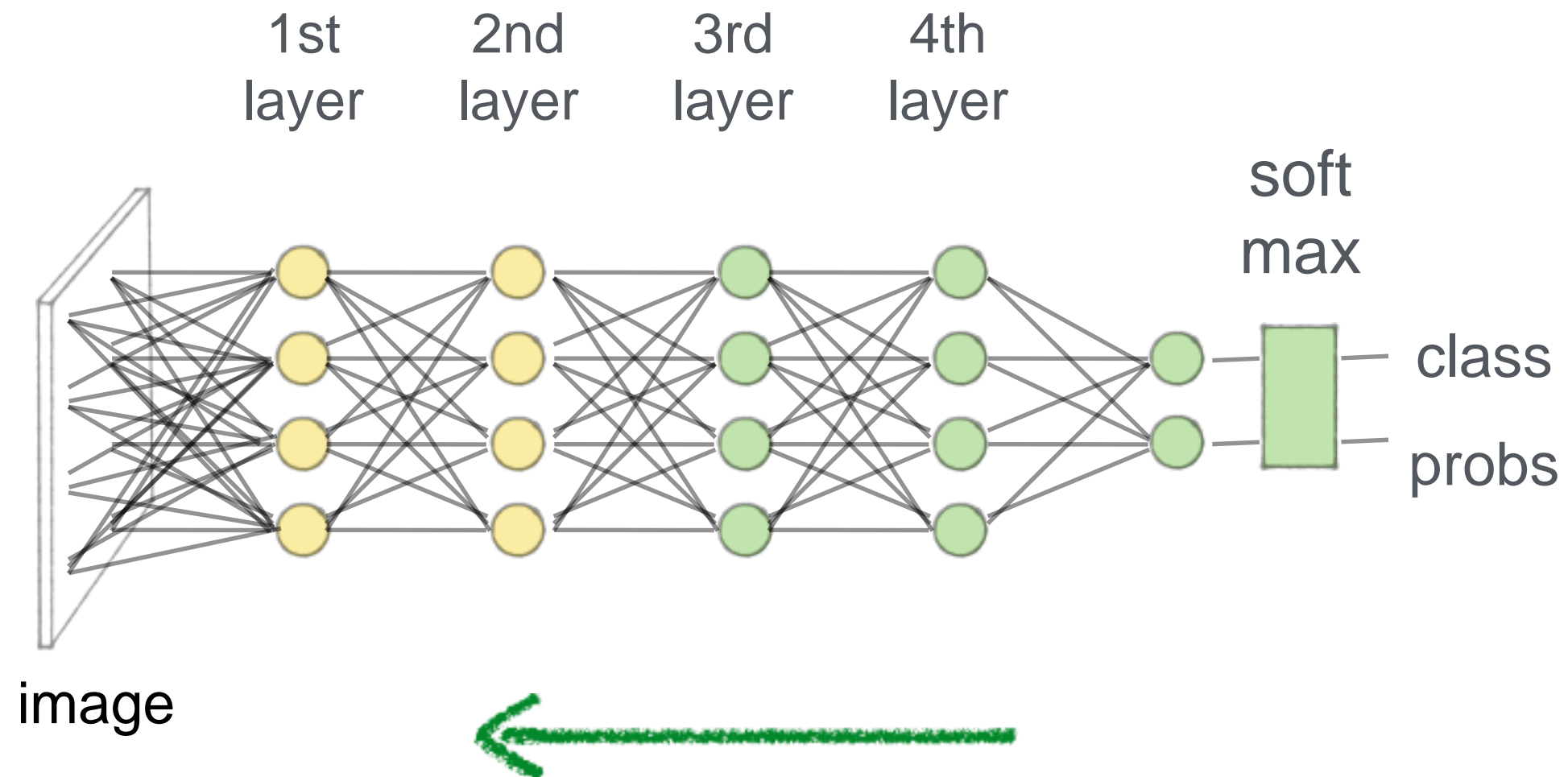
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

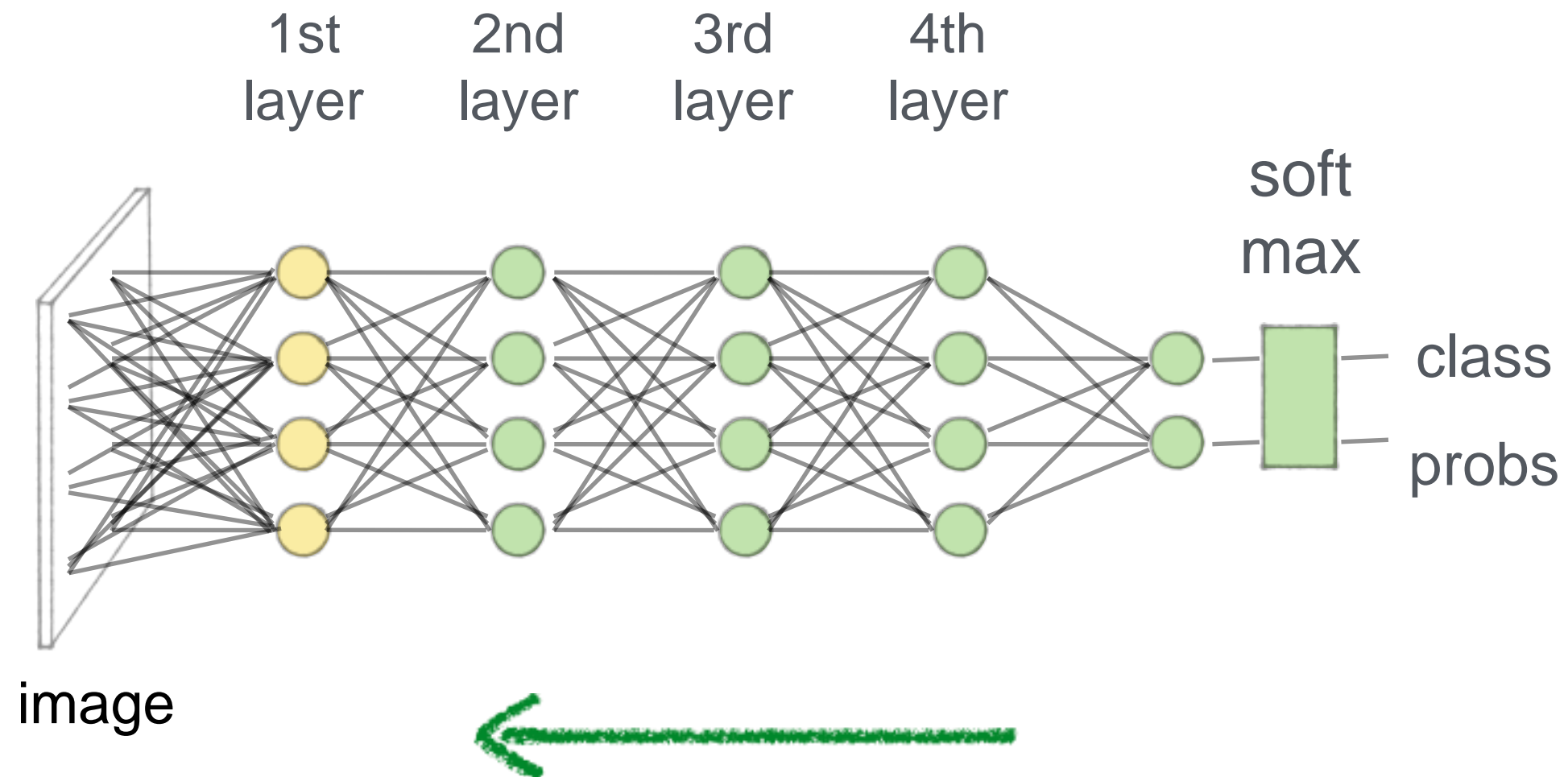
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

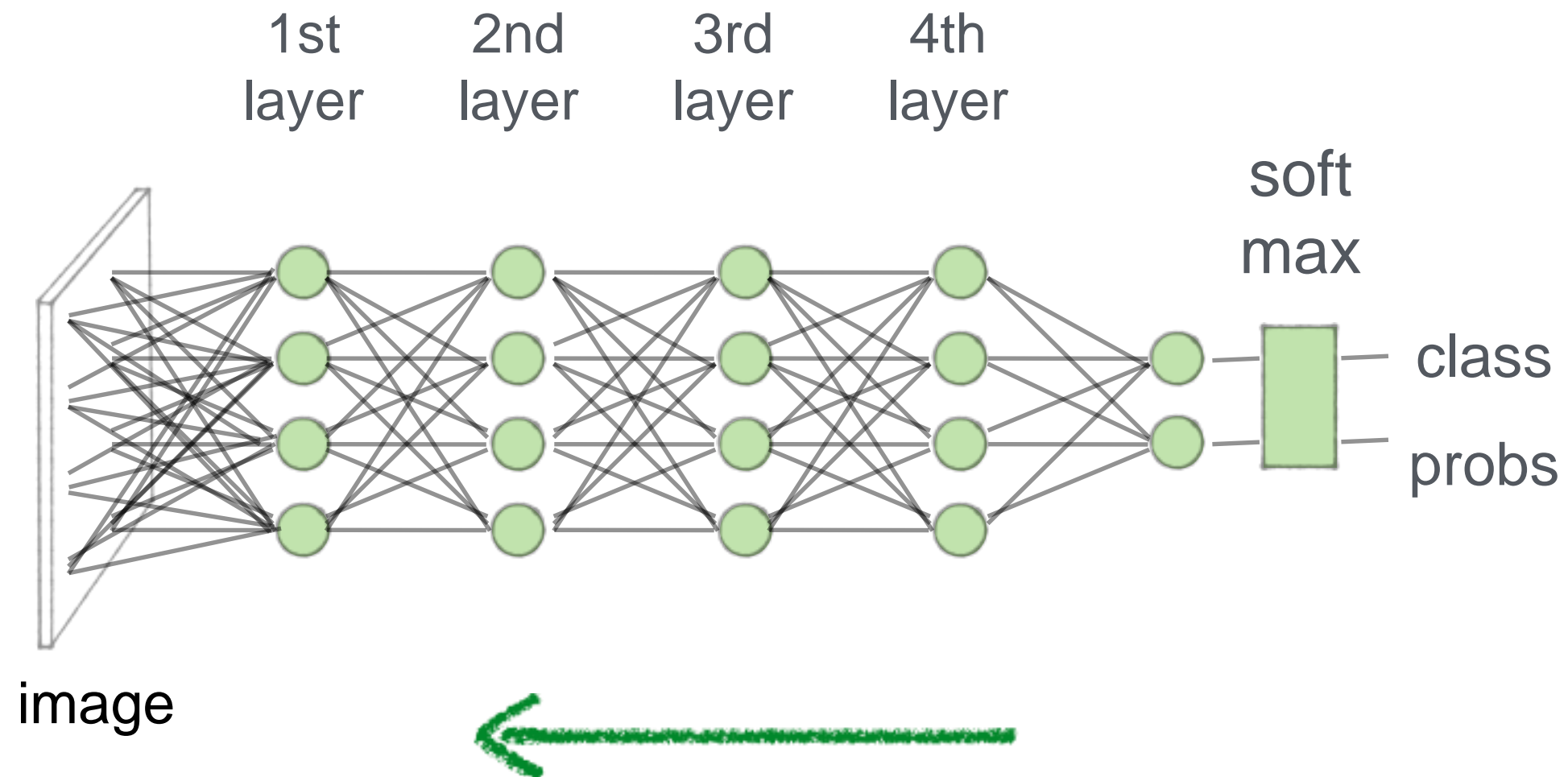
Last Lecture



Backpropagation (backprop)

$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

Last Lecture



Backpropagation (backprop)

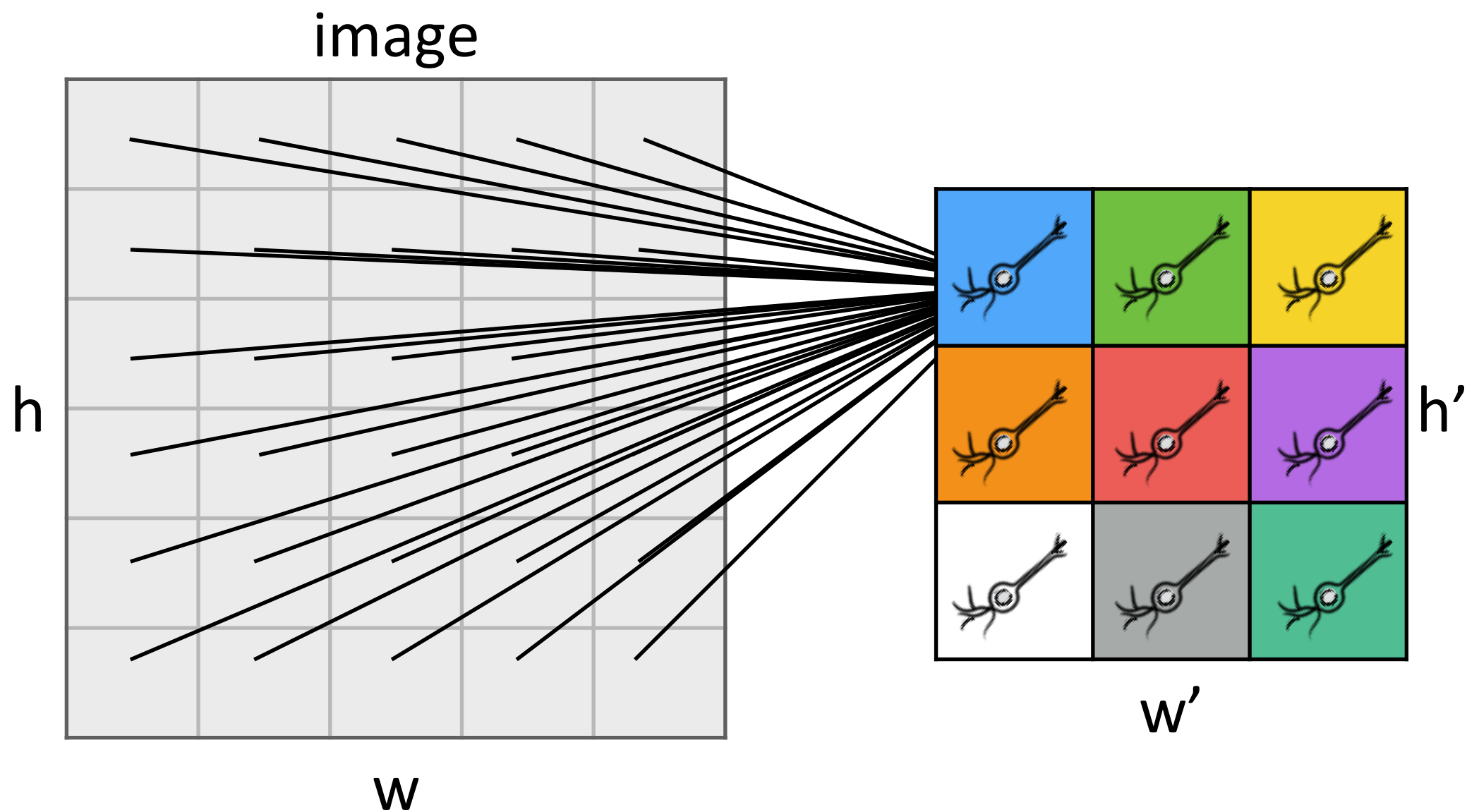
$$\frac{\partial L_i}{\partial w_k} = \underbrace{x_k}_{\text{input}} \underbrace{\xi'}_{\text{error}} \sum_j \tilde{w}_j \underbrace{\frac{\partial L_i}{\partial \tilde{y}_j}}_{\text{error}}$$

Today

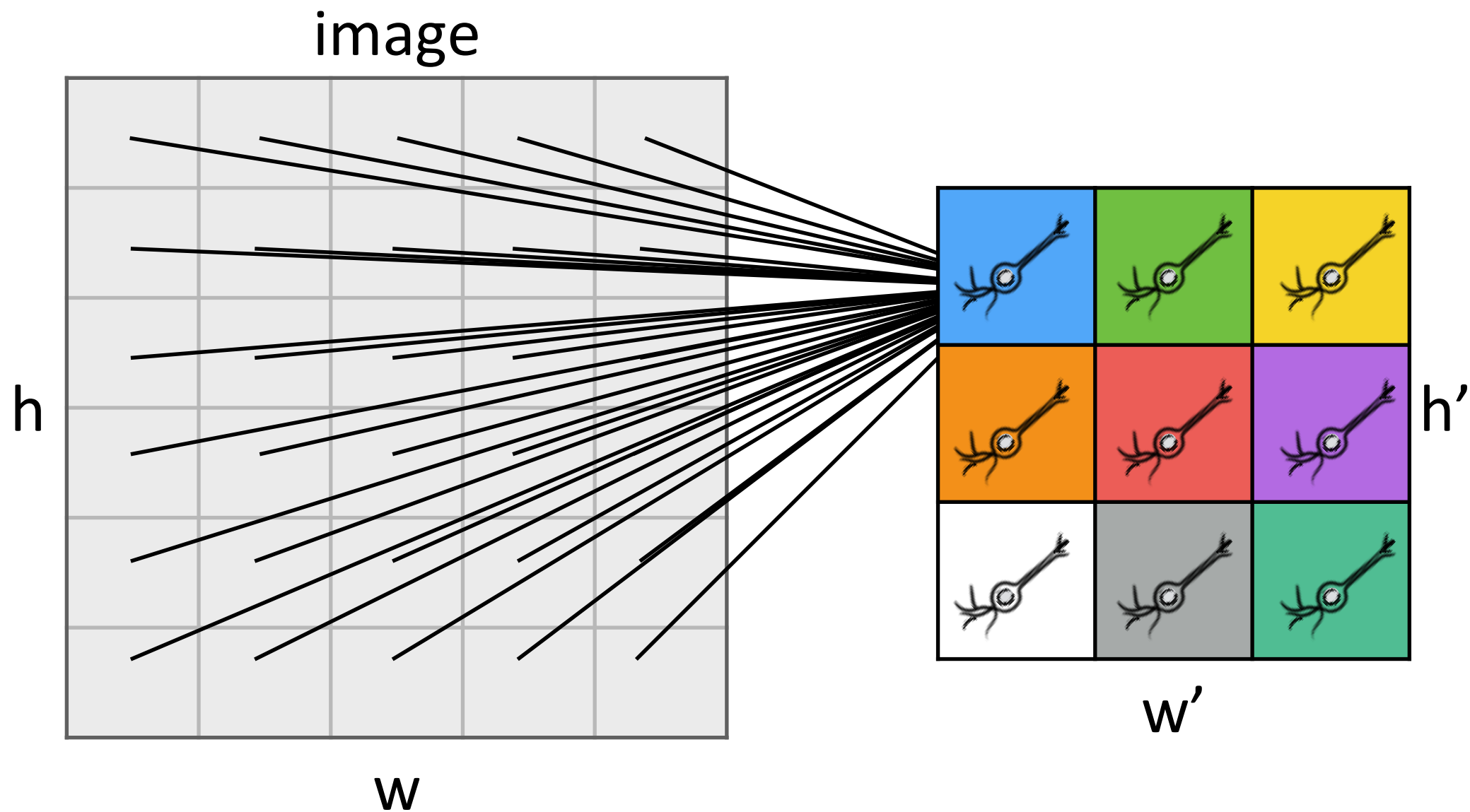
- Recap: Convolutional Neural Networks
- Fully Convolutional Neural Networks
- Convolutional Layers
- Overfitting, Part II

Convolutional Neural Networks

Fully Connected Layers

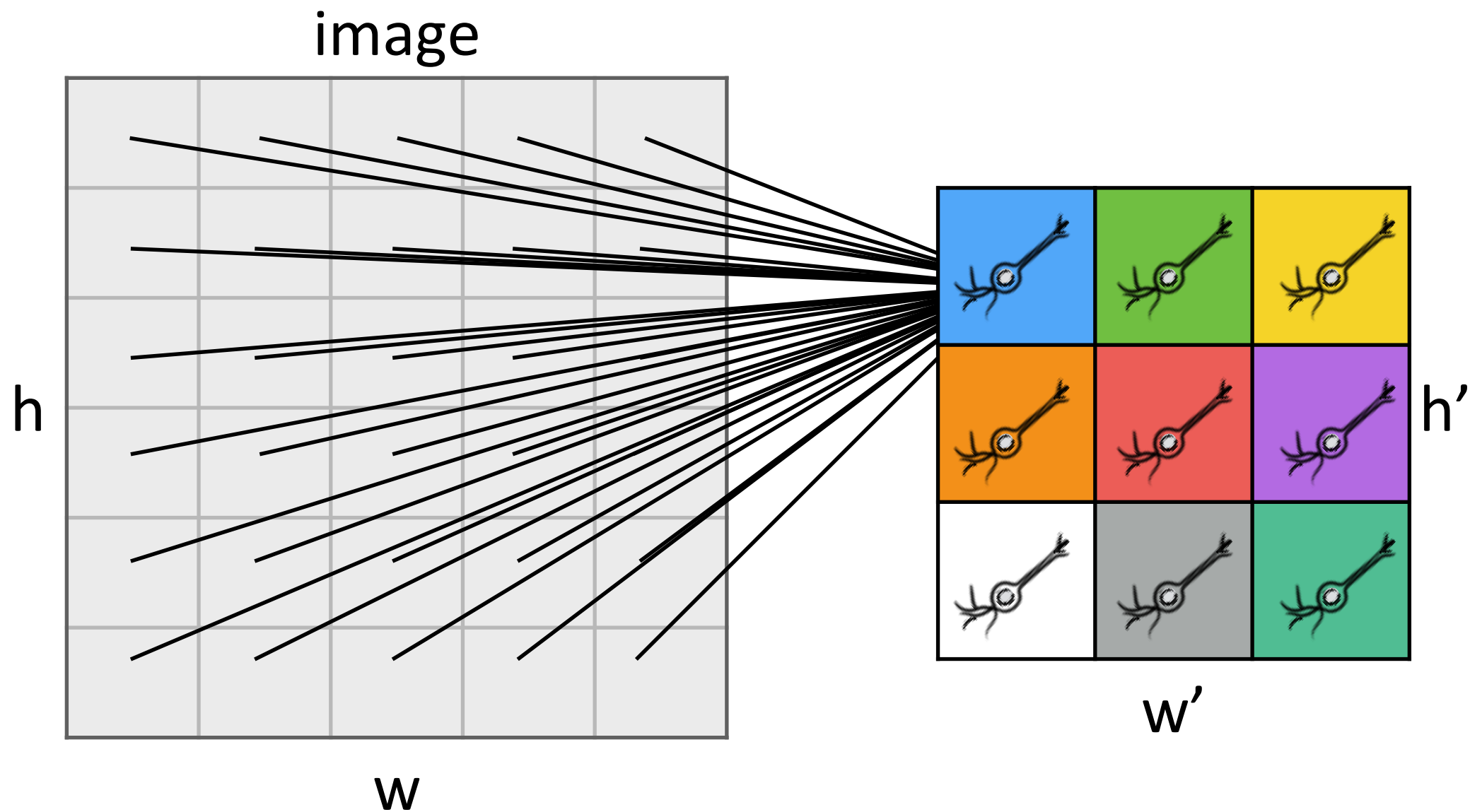


Fully Connected Layers



Number parameters: $w \cdot h \cdot w' \cdot h'$

Fully Connected Layers

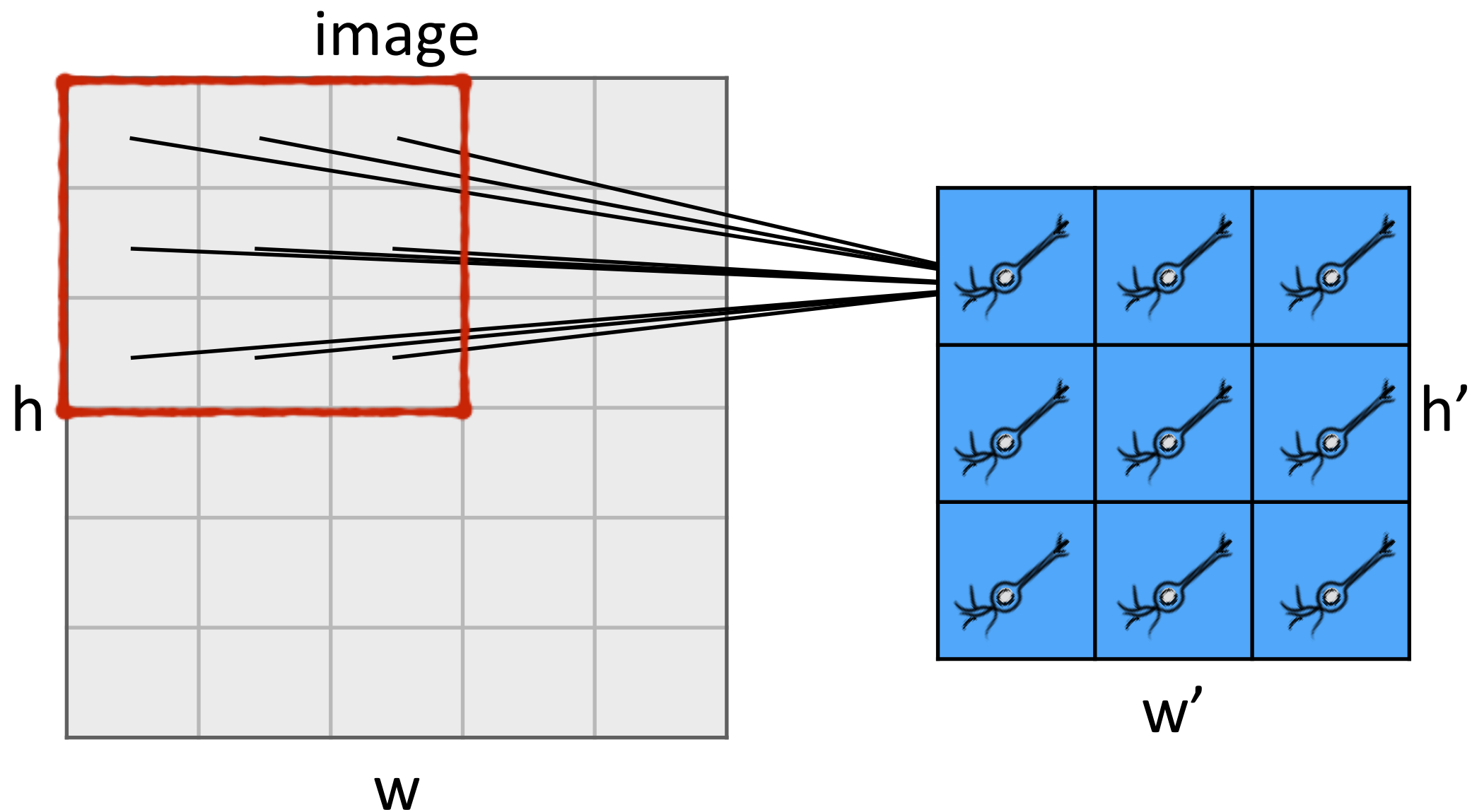


Redundancy

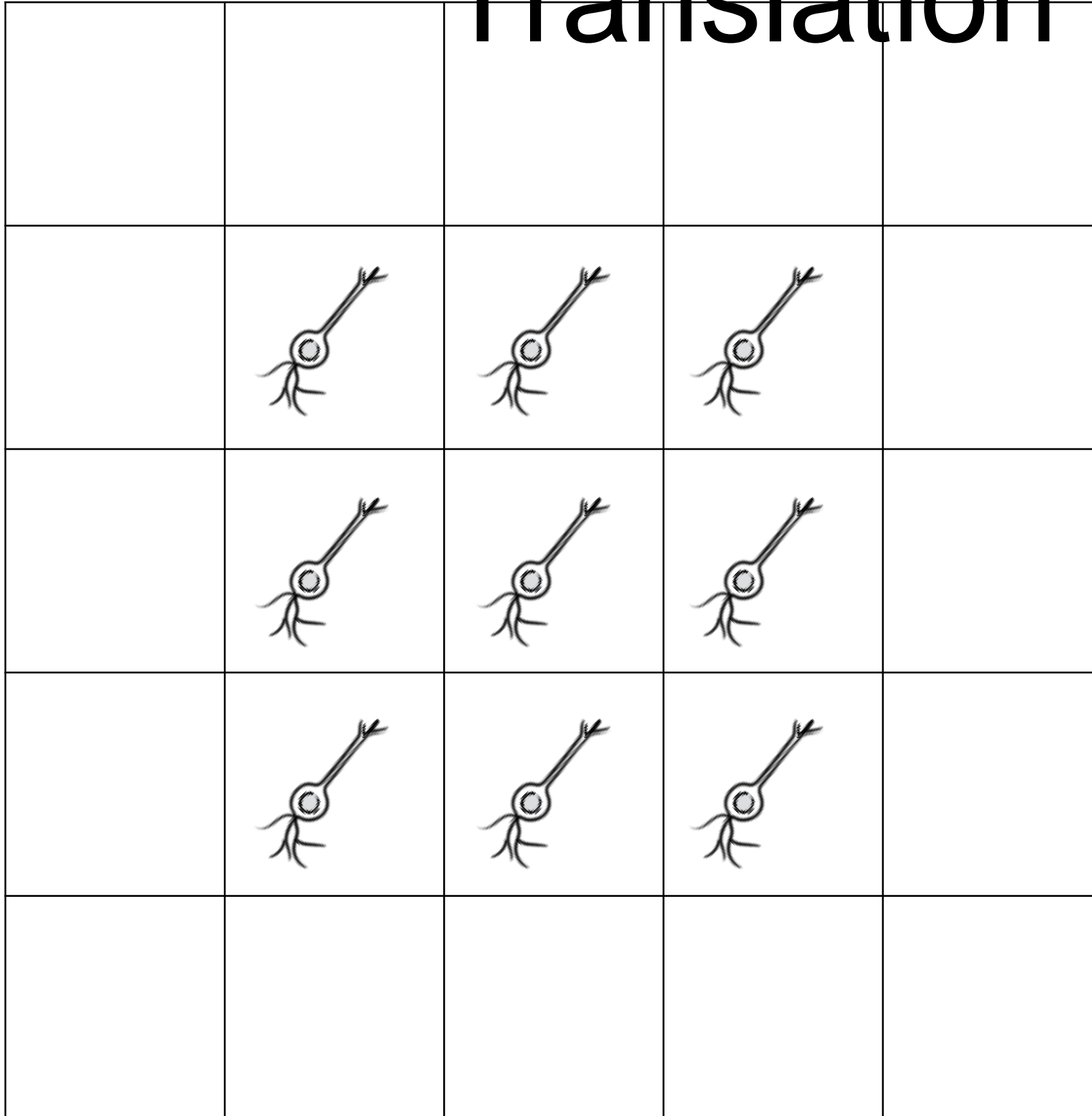
Number parameters: $w \cdot h \cdot w' \cdot h'$

But no locality

Convolutional Layers

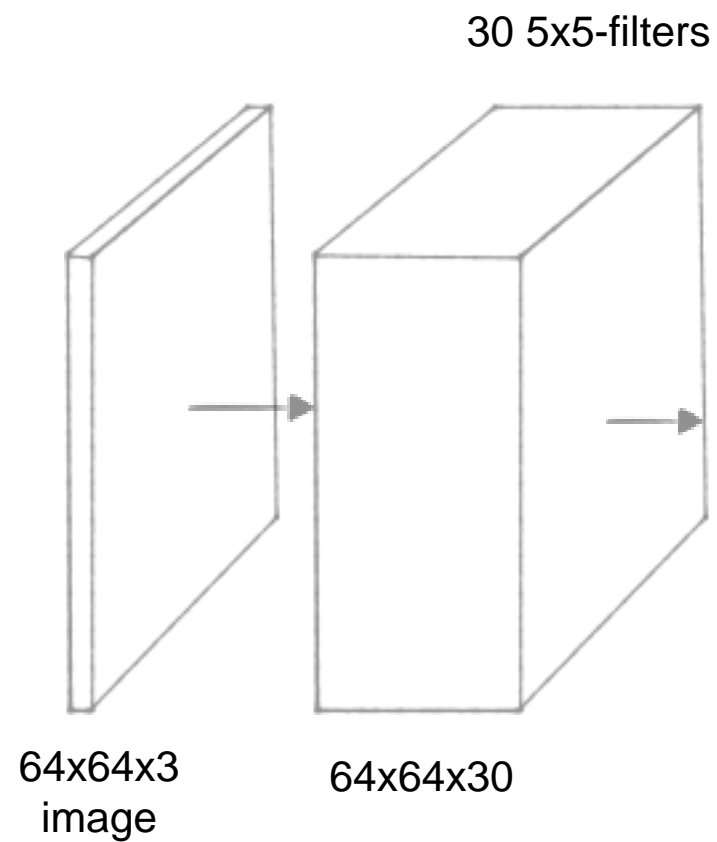


Translation Invariance



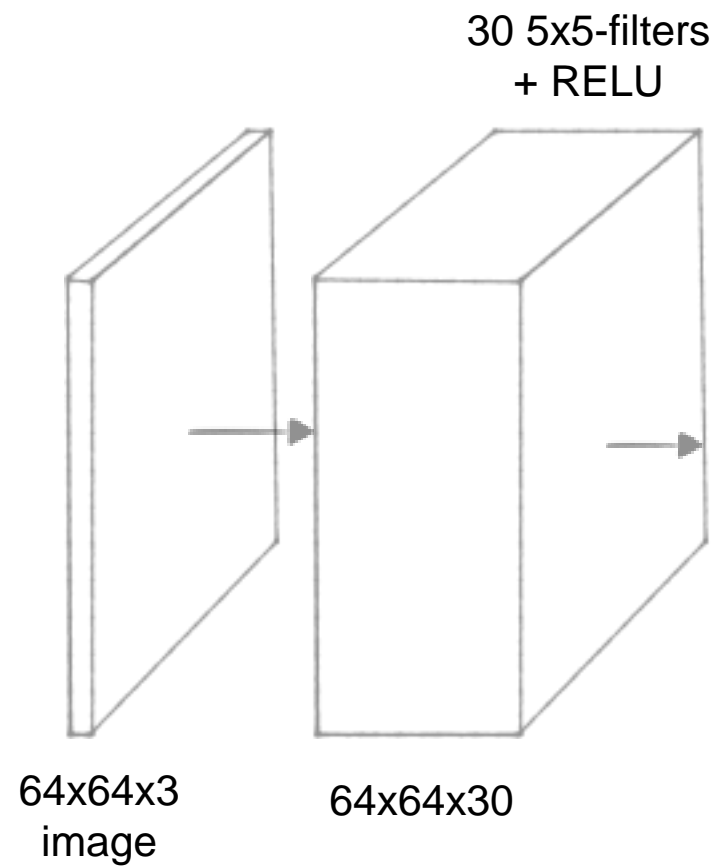
$$I \star \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

Network Structure



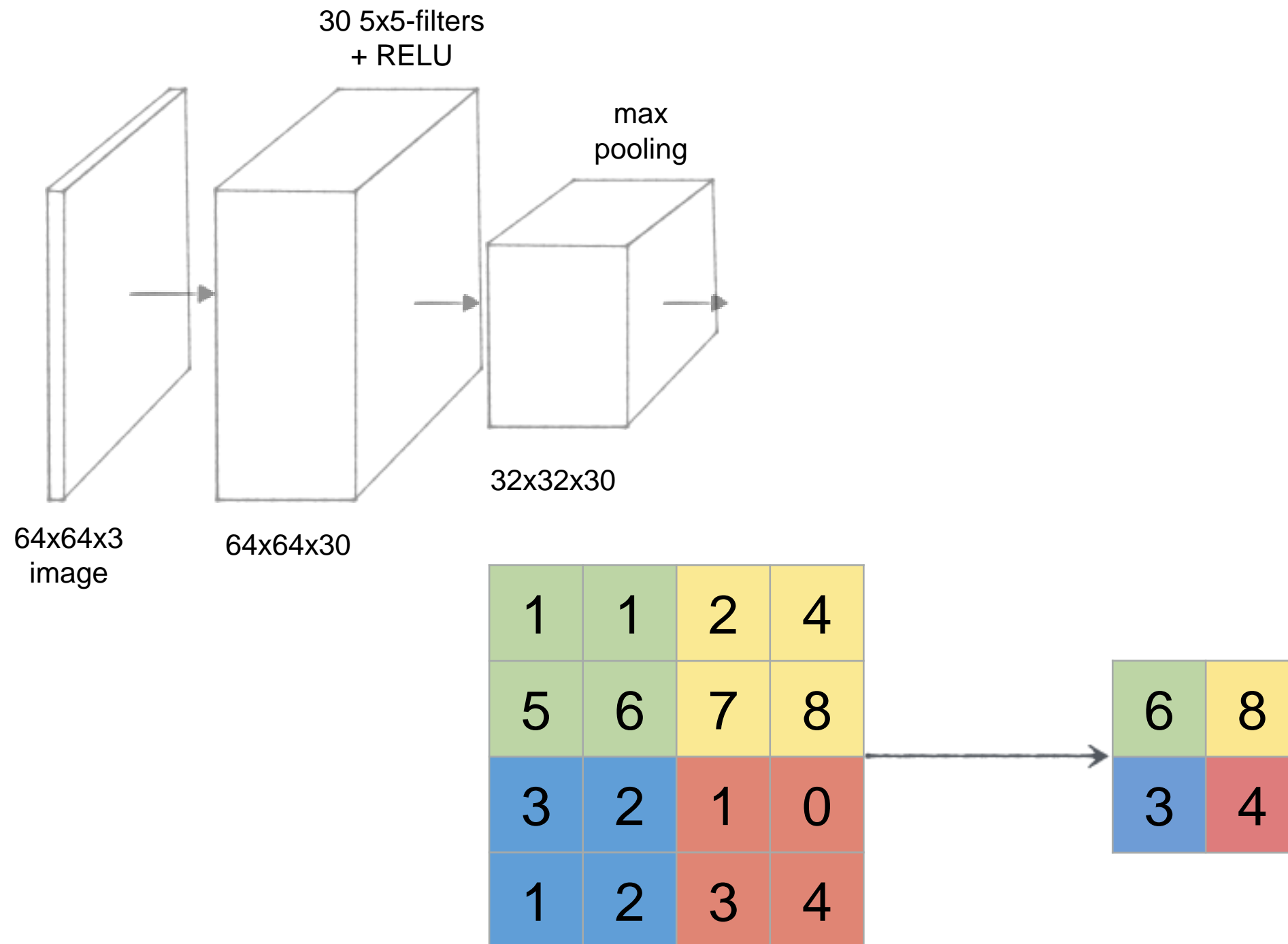
$$I \star w$$

Network Structure

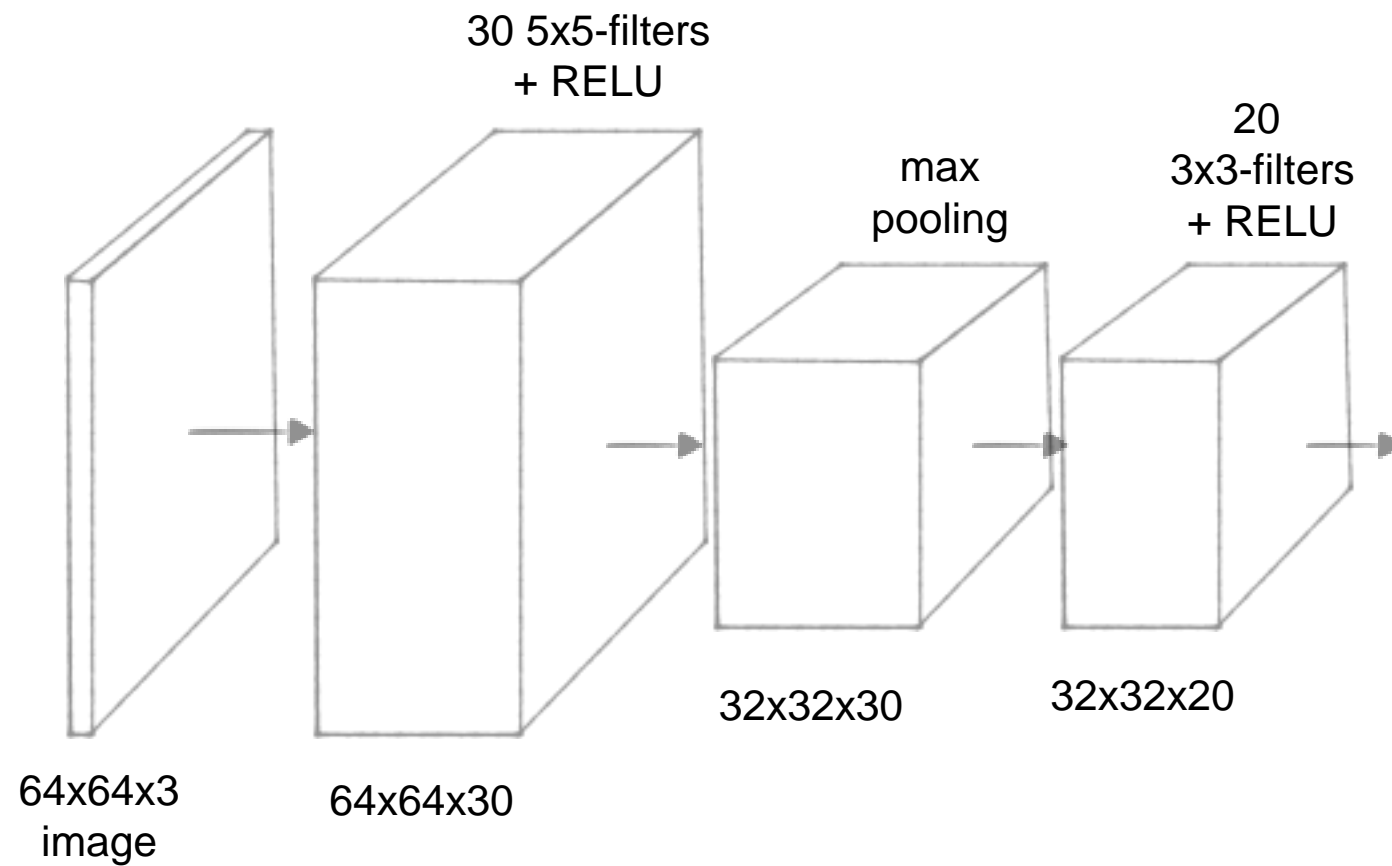


$$\max\{I \star w, 0\}$$

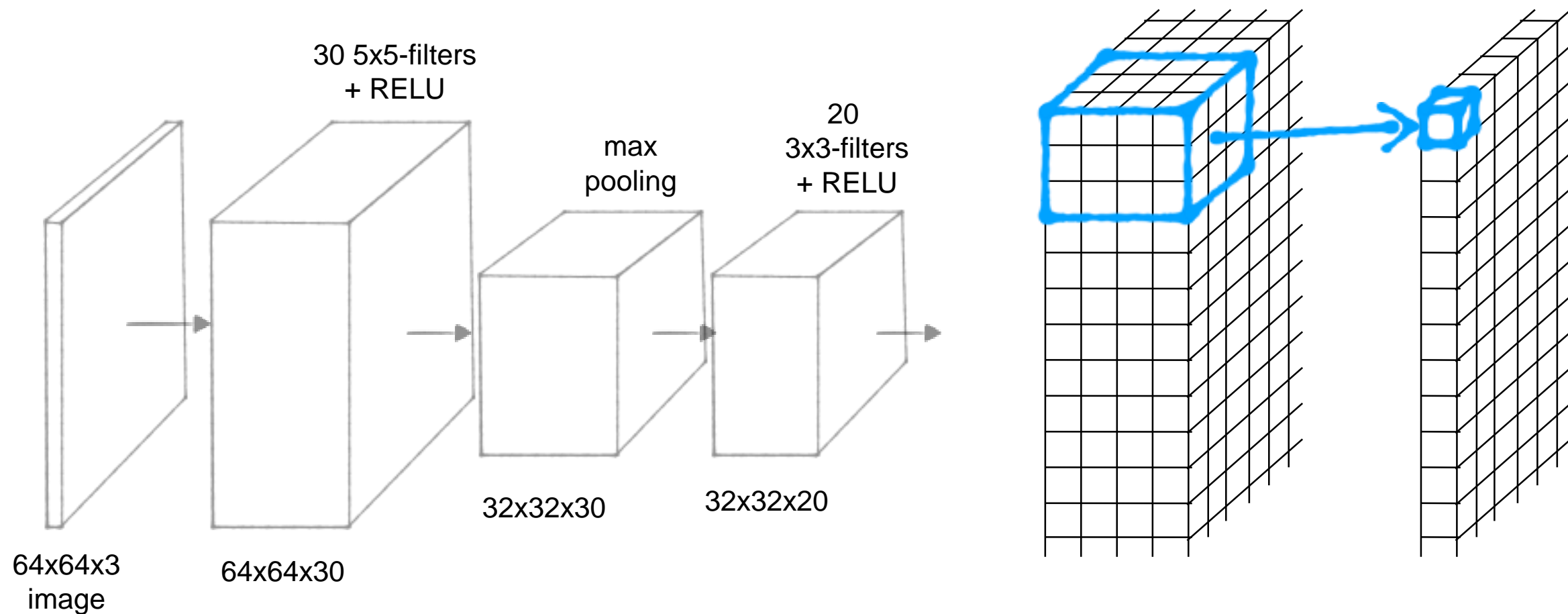
Network Structure



Network Structure



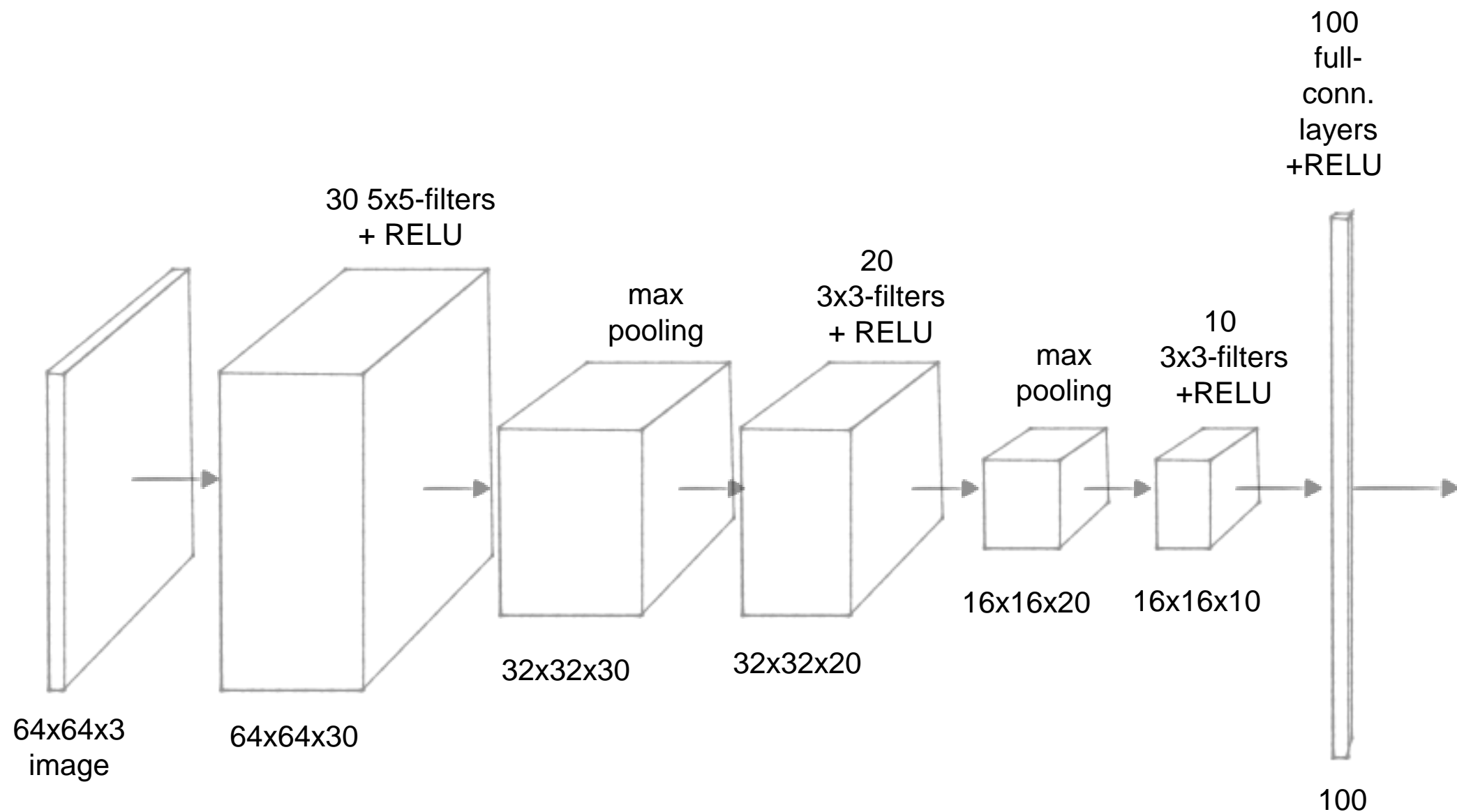
Network Structure



implemented as matrix-vector multiplication:

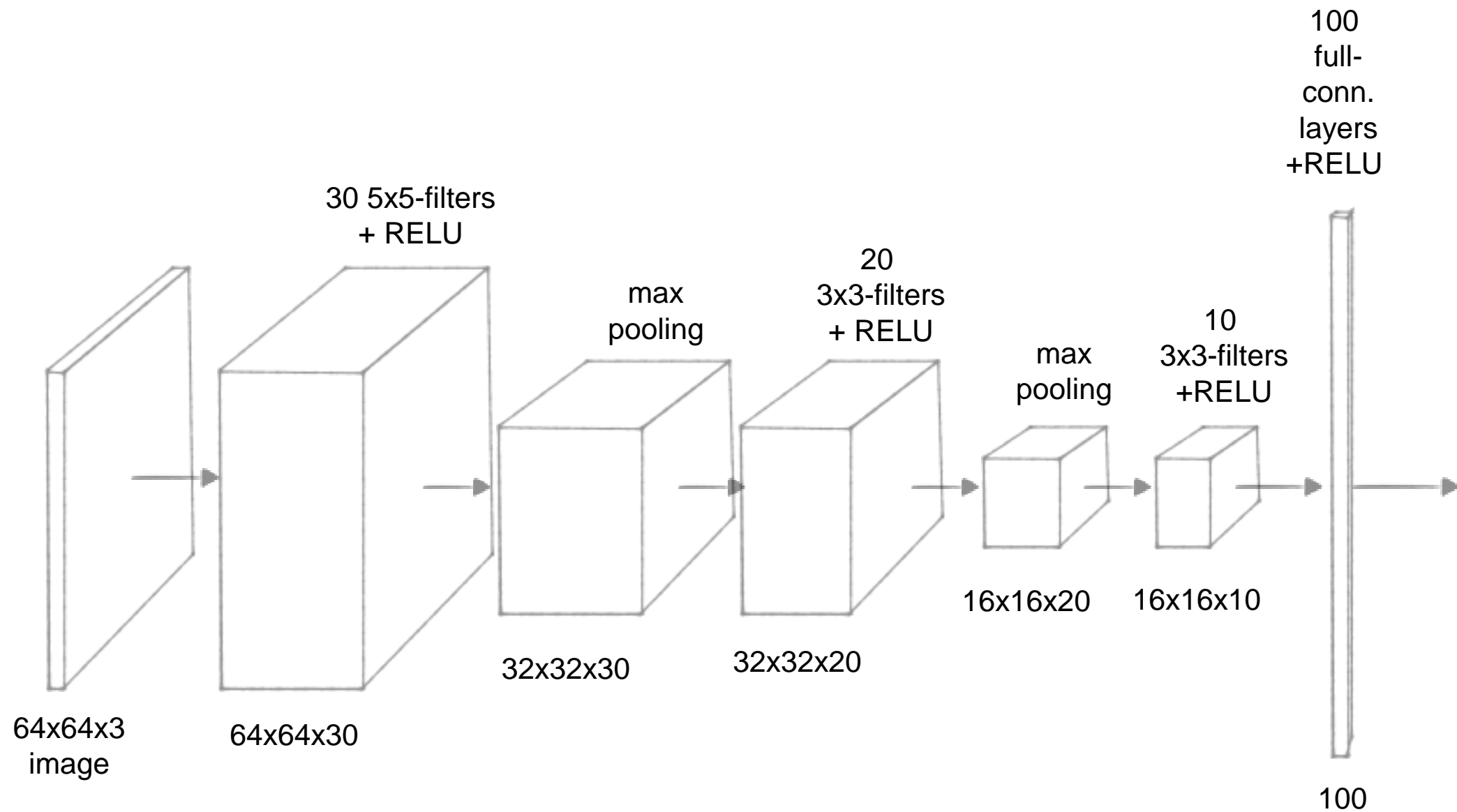
$$\text{vec}(\mathcal{J}) = F \text{vec}(\mathcal{I})$$

Network Structure



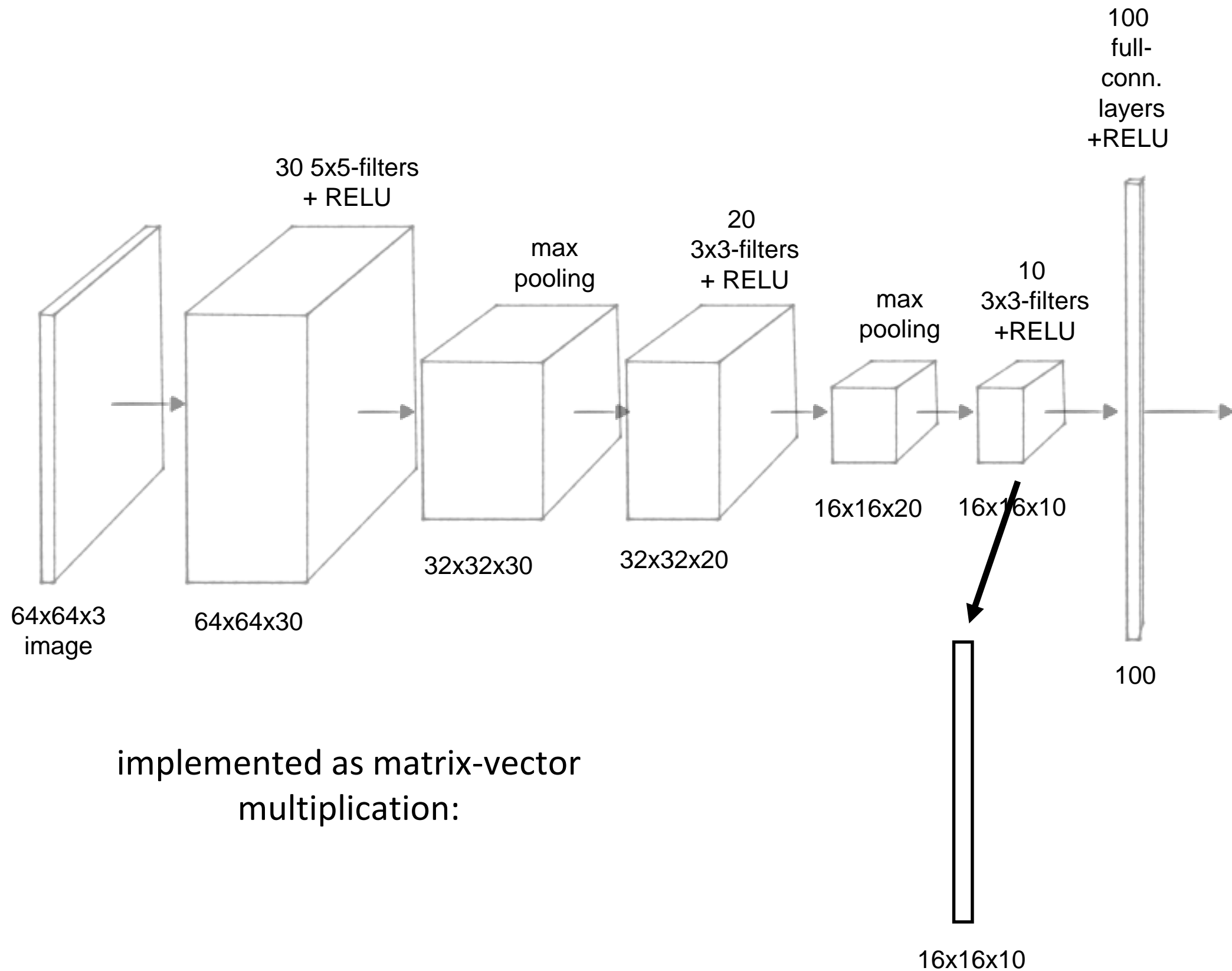
$$\sum_{u=1}^{16} \sum_{v=1}^{16} \sum_{i=1}^{10} w_{uvi} x_{uvi}$$

Network Structure

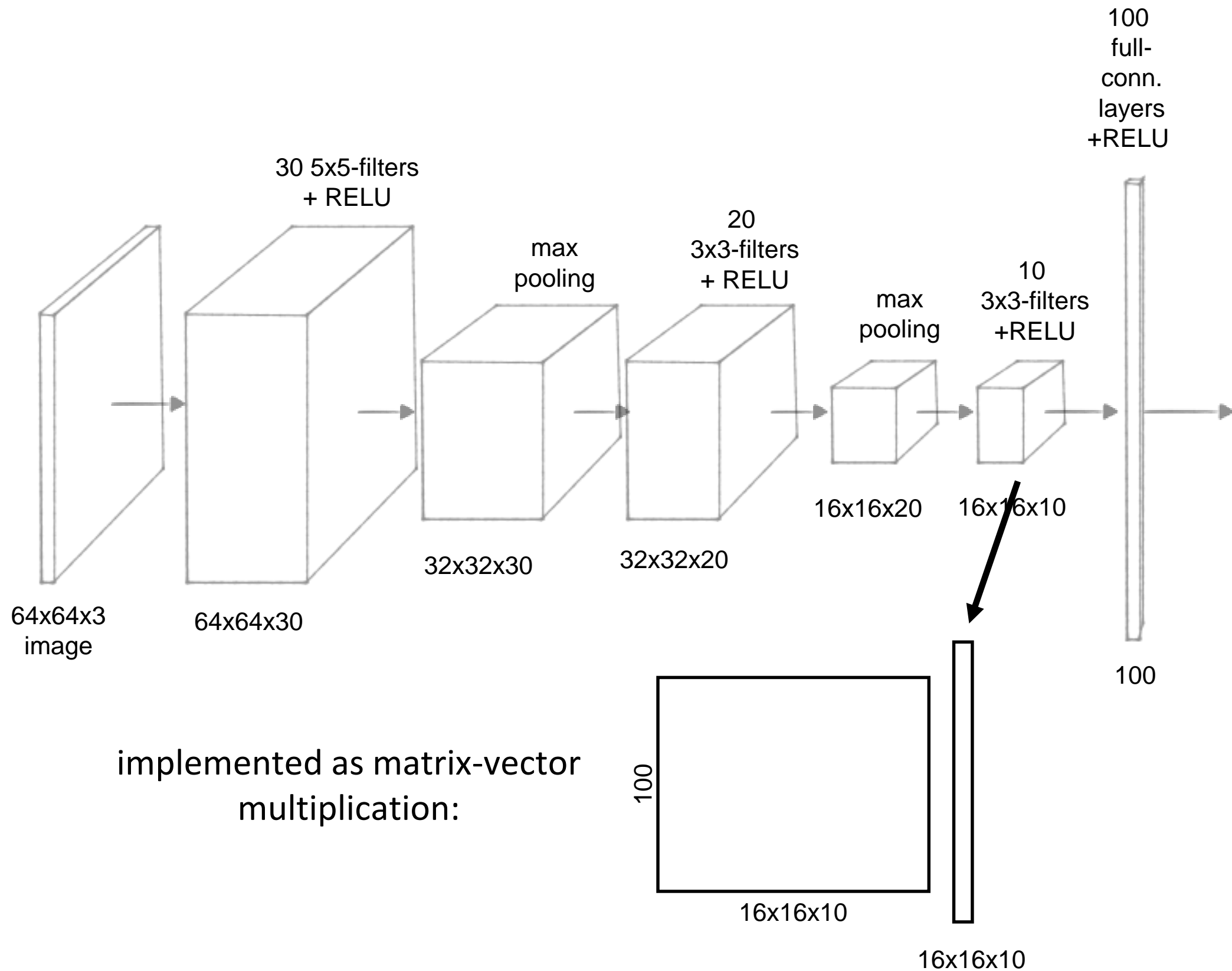


implemented as matrix-vector
multiplication:

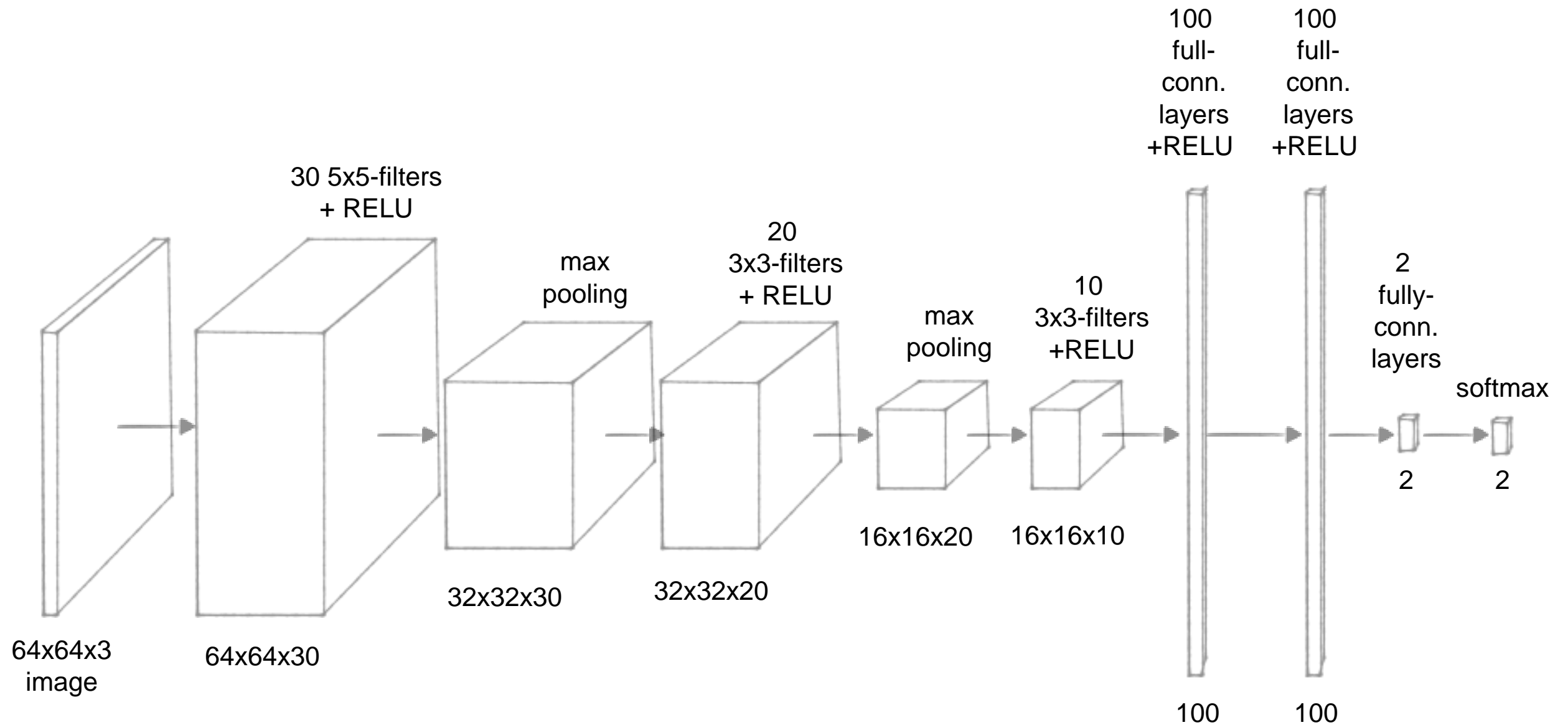
Network Structure



Network Structure



Network Structure



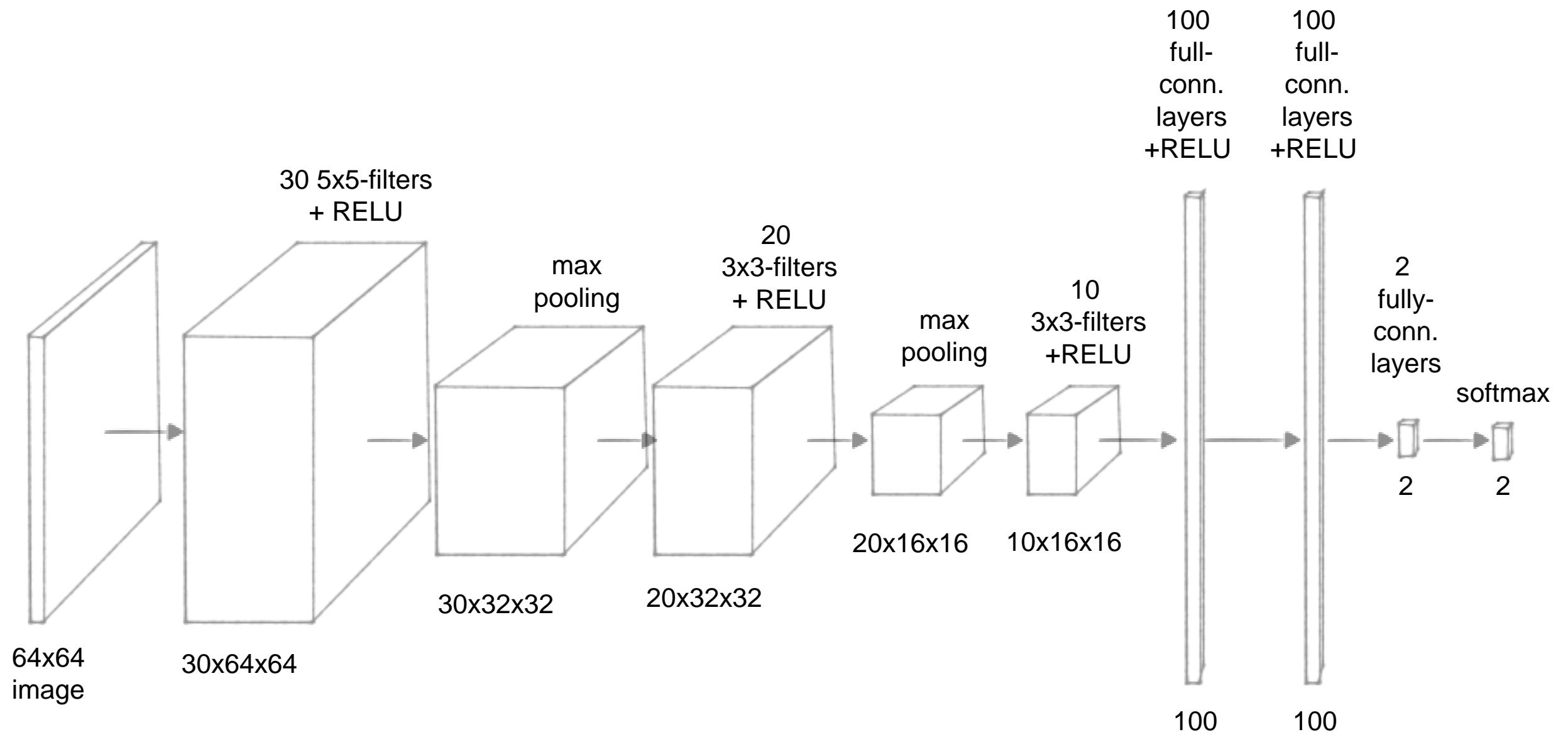
Fully Convolutional Neural Networks

Recognition / Classification

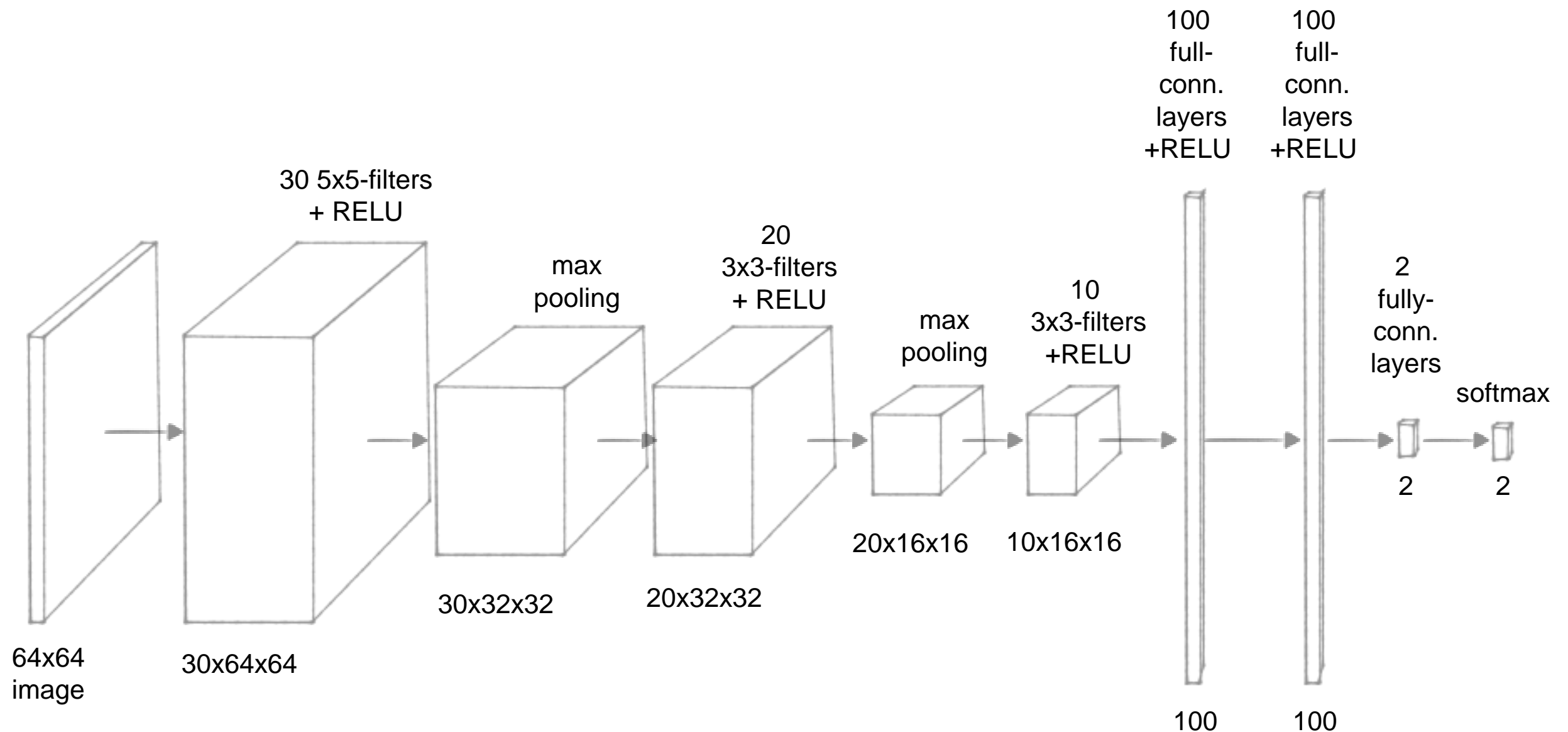


Is there a robin in this image?

Convolutional Neural Network for Classification

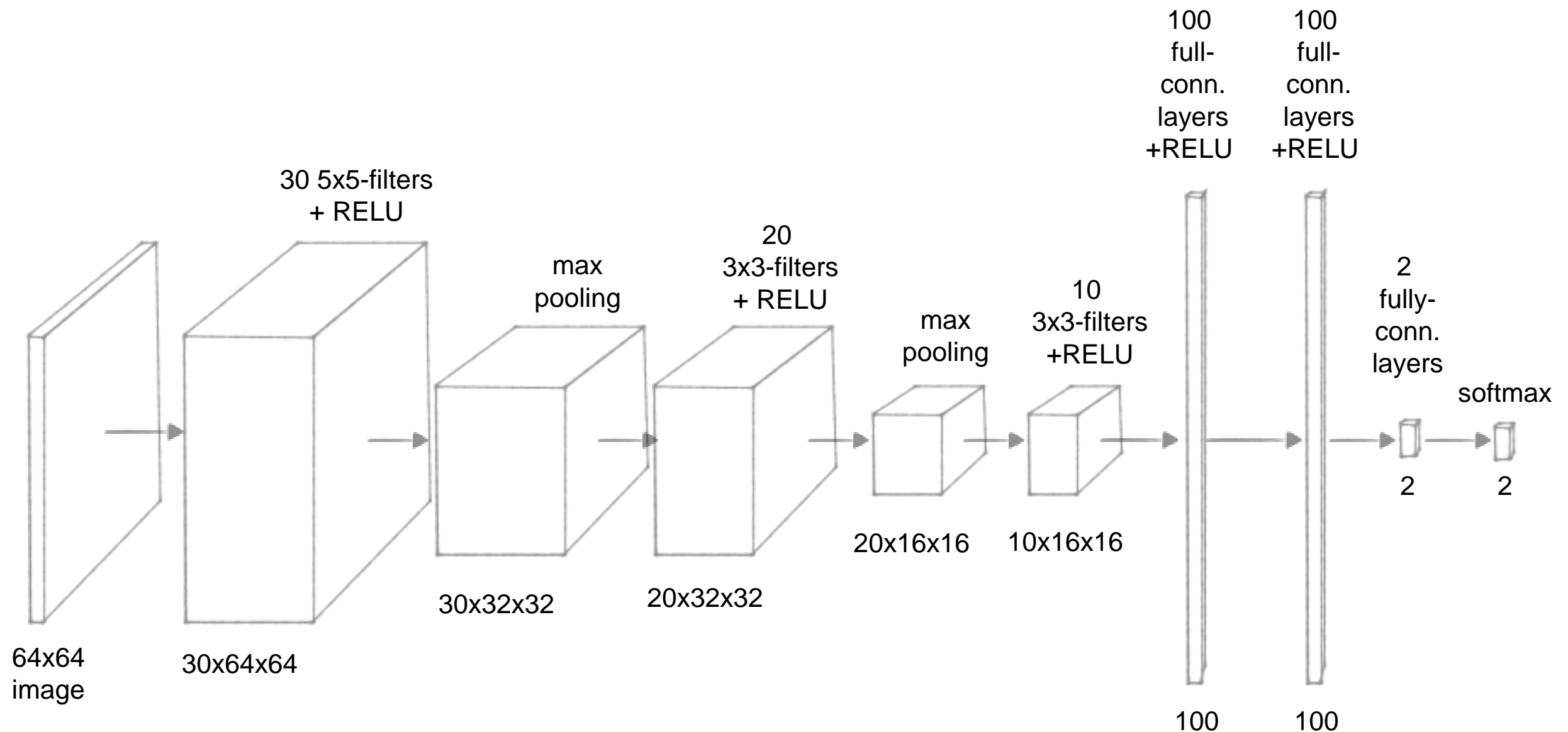


Convolutional Neural Network for Classification



Input image size needs to be fixed during training!

Convolutional Neural Network for Classification



Input image size needs to be fixed during training!

How to get a per-pixel decision?

Instance Segmentation

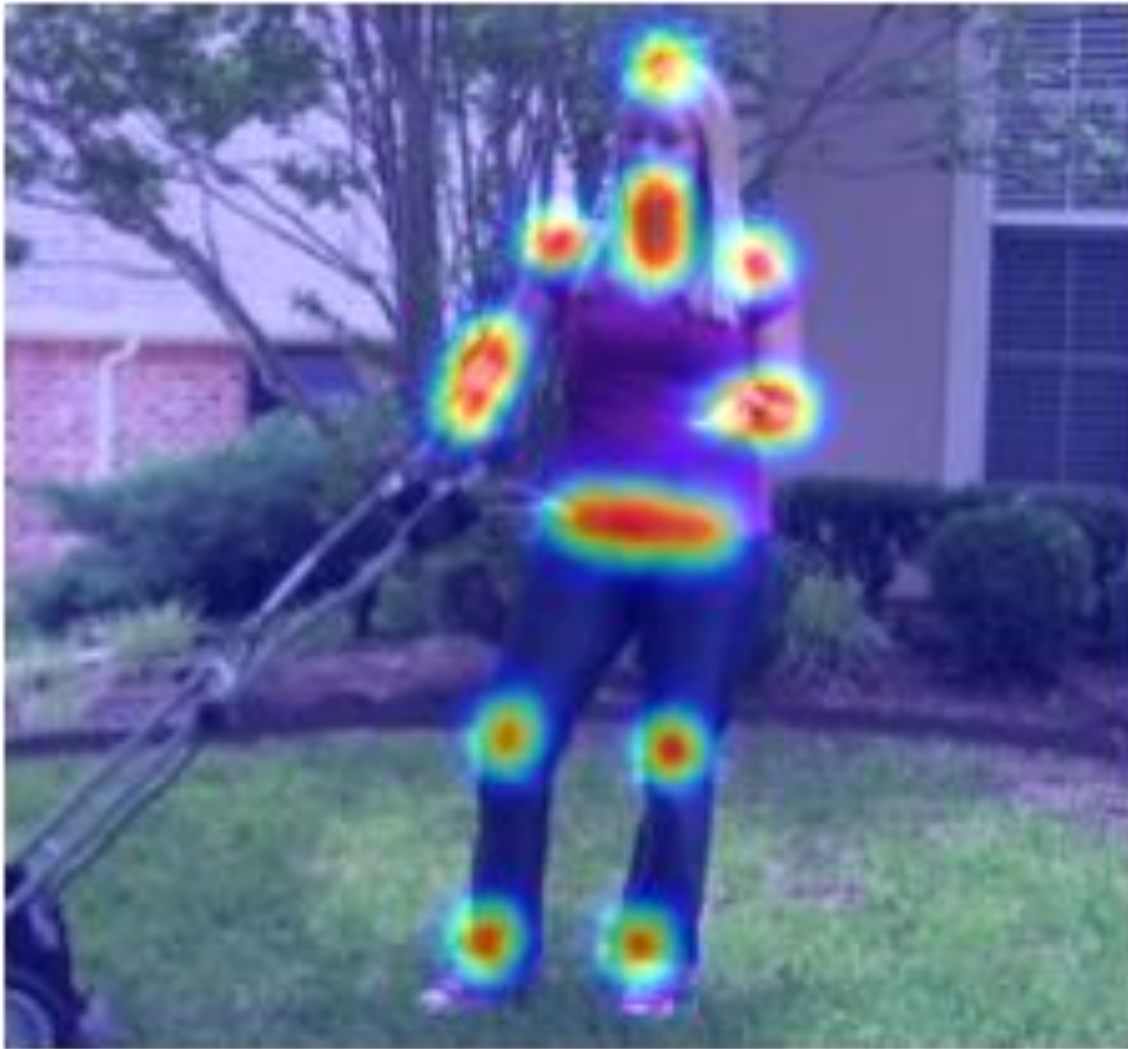


Semantic Segmentation



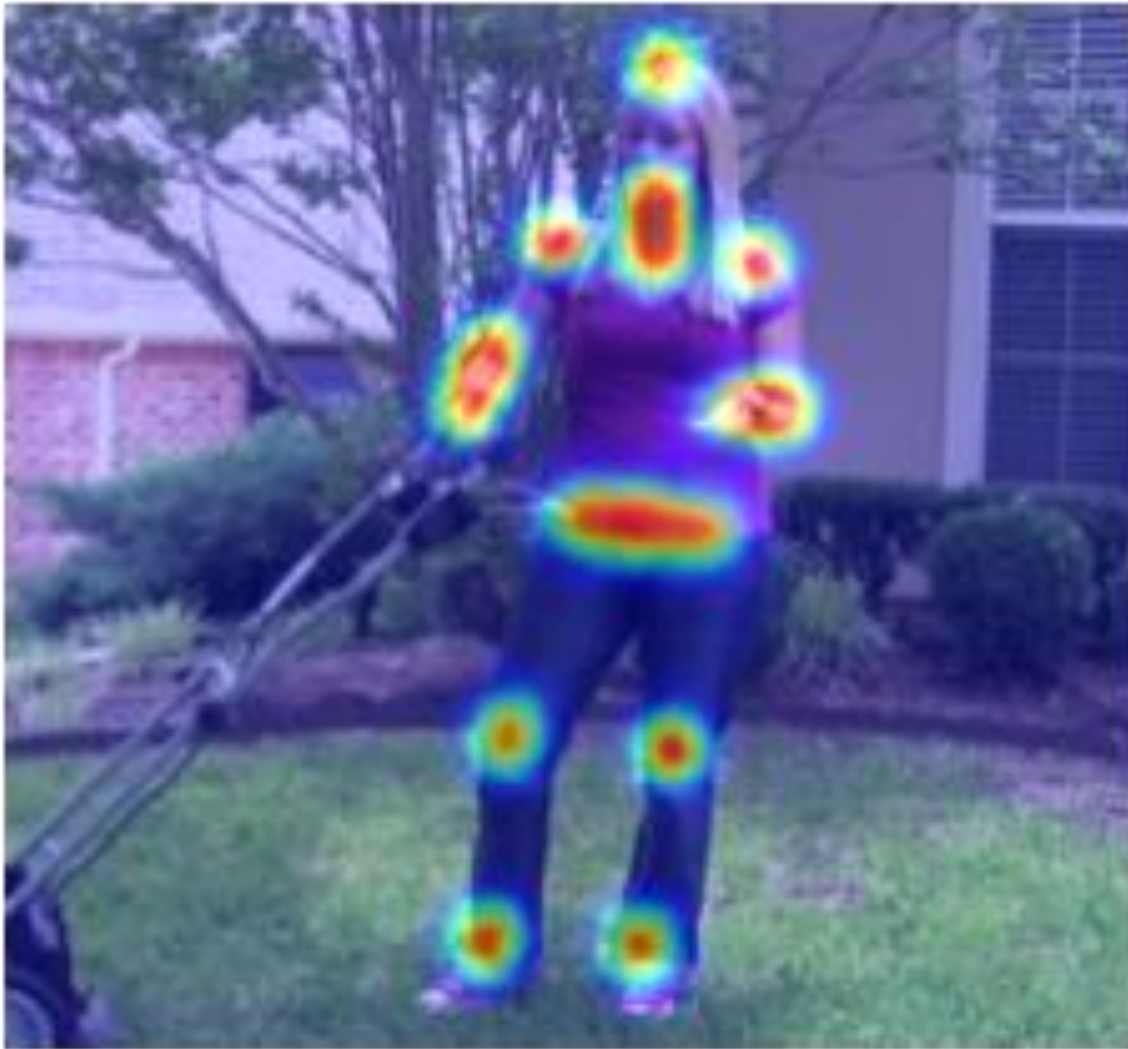
image source: [Mapillary](https://www.mapillary.com)

Human Pose Estimation



Needs multiple decisions.

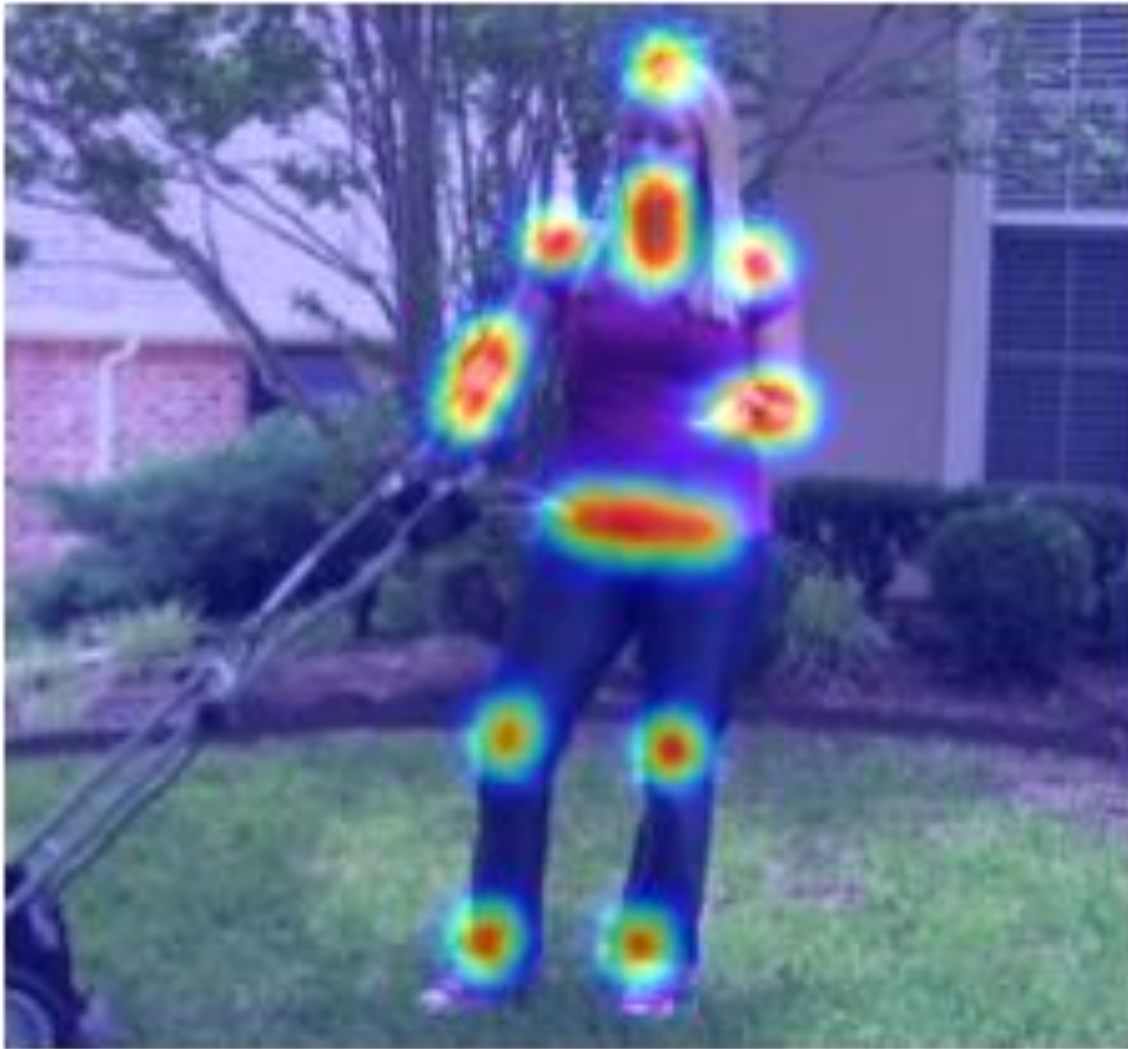
Human Pose Estimation



Needs multiple decisions.

Just apply network in sliding window fashion?

Human Pose Estimation



Needs multiple decisions.

Just apply network in sliding window fashion? Works, but slow!

Fully-Convolutional Networks

Faster approach:



From Fully-Connected to Fully-Convolutional

- Fully-connected layer is still a convolutional layer
- Can convert a network with fully-connected into fully-convolutional network:
 - Simply reshape matrix from fully-connected layer into a filter
 - Produces the same result on the input volume of the fully-connected layer
 - See example on black board

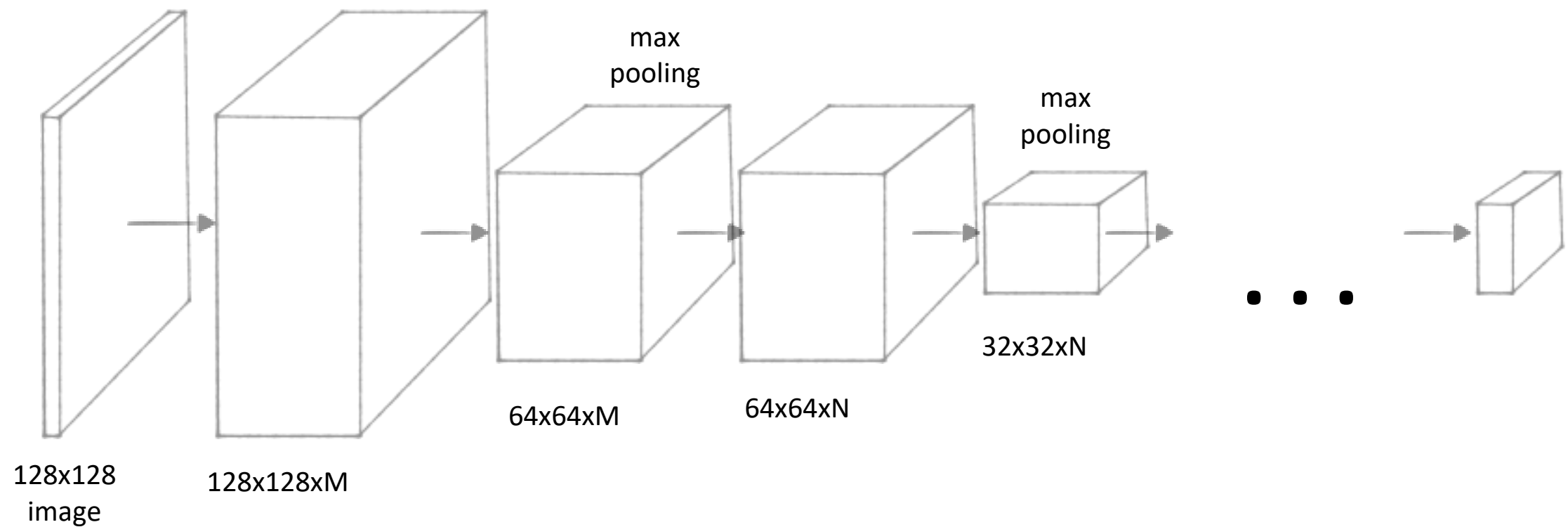
Fully-Convolutional Networks

Faster approach:

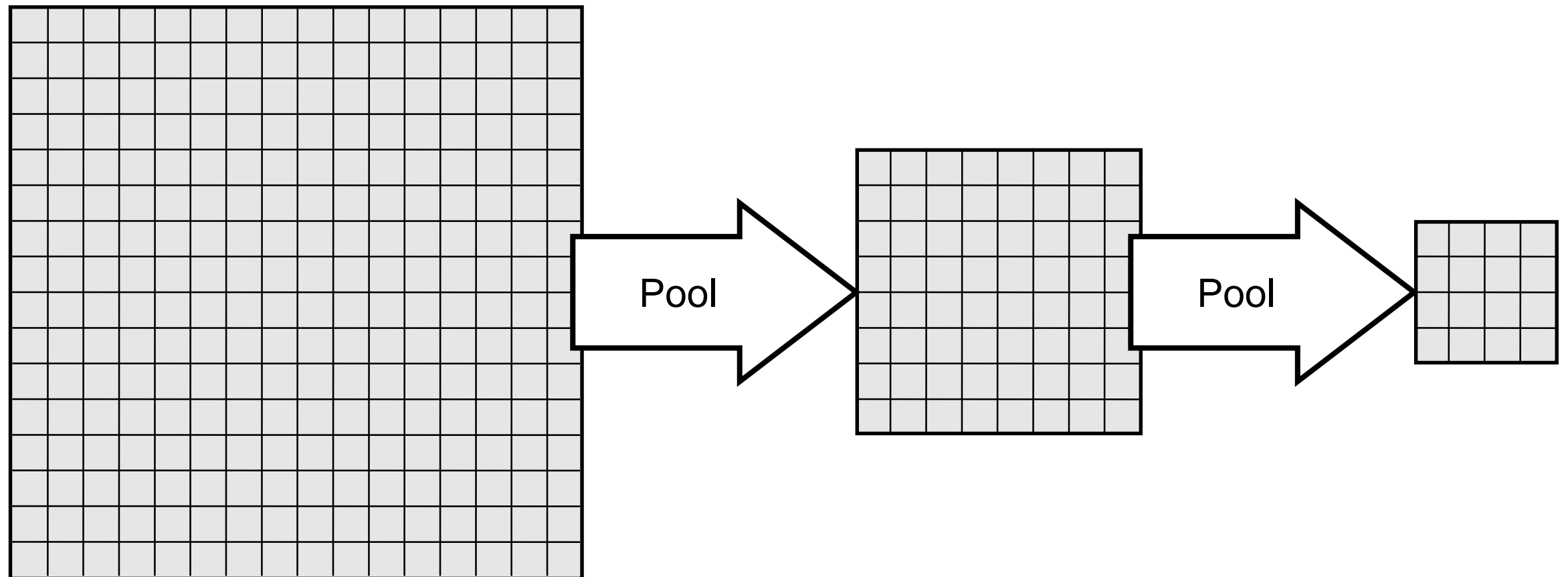


But: Memory intensive at full resolution

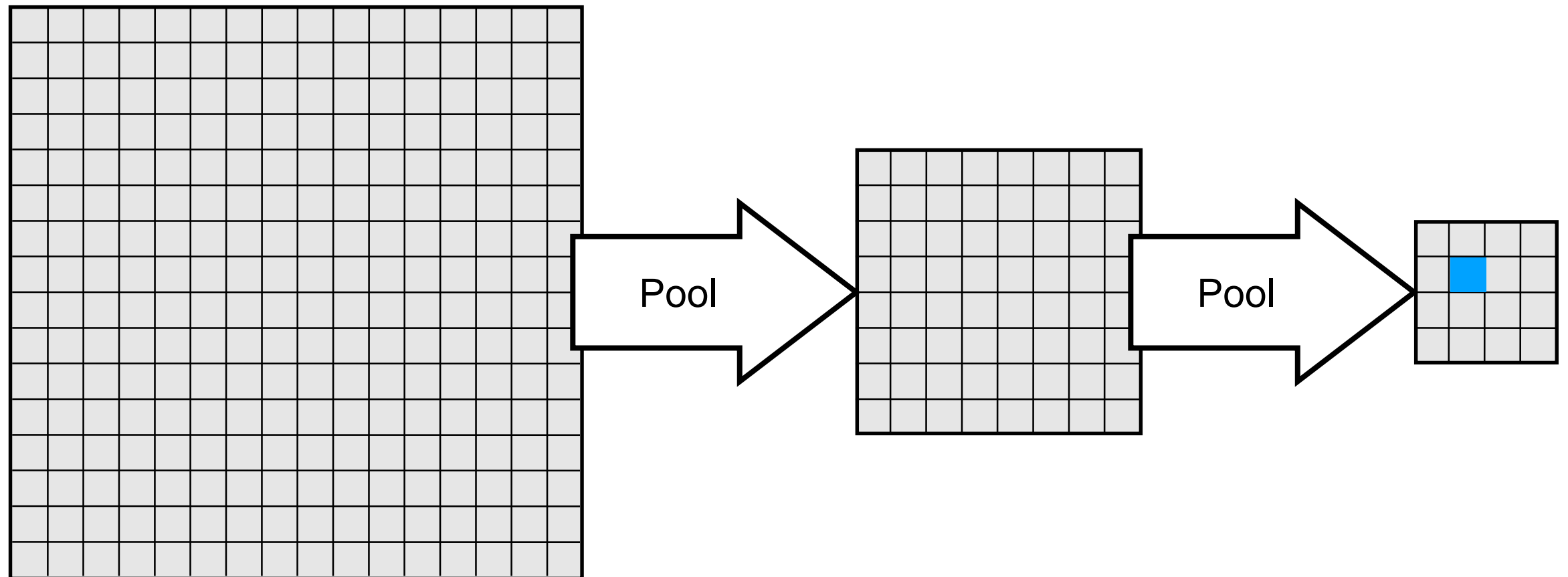
Fully-Convolutional Network for Detection



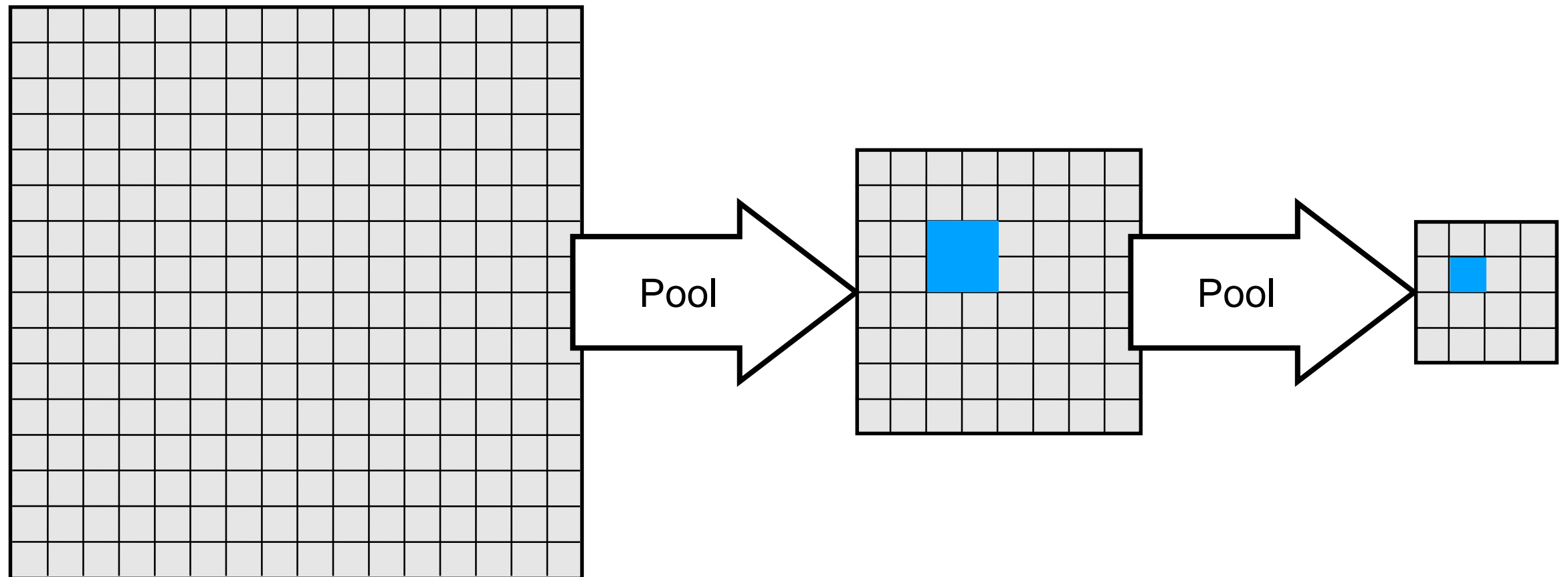
Further Benefit of (Max-)Pooling



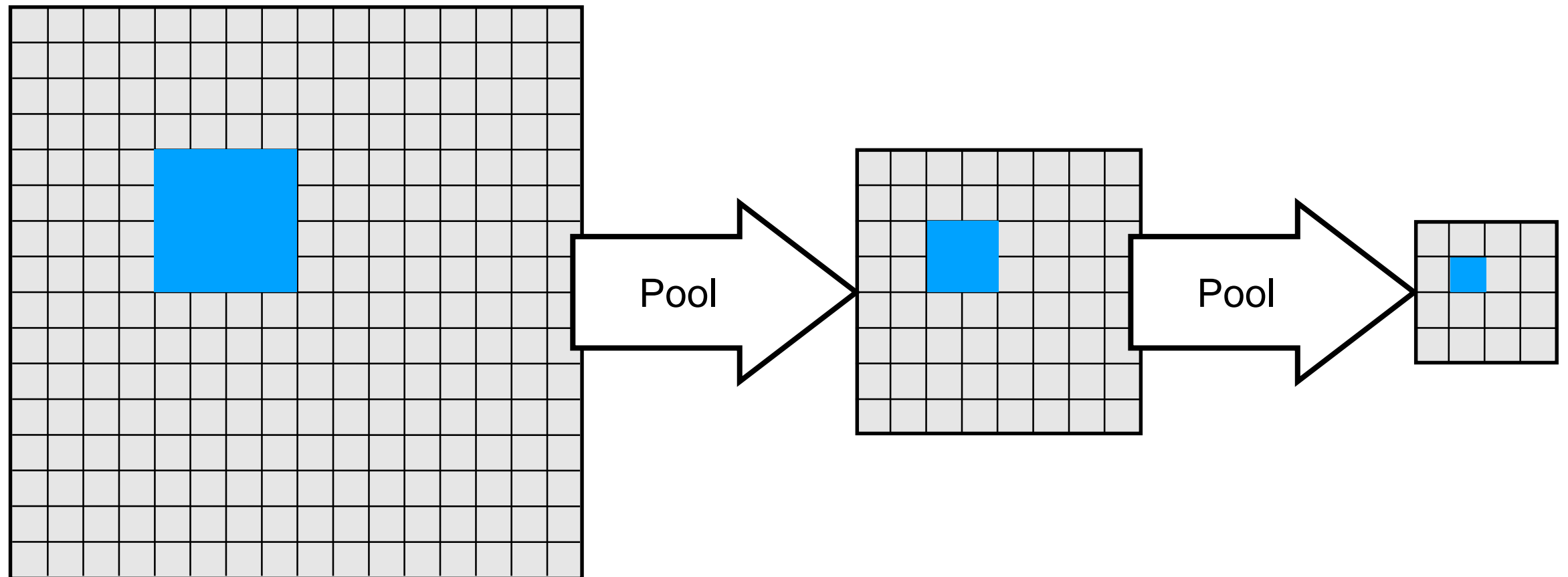
Further Benefit of (Max-)Pooling



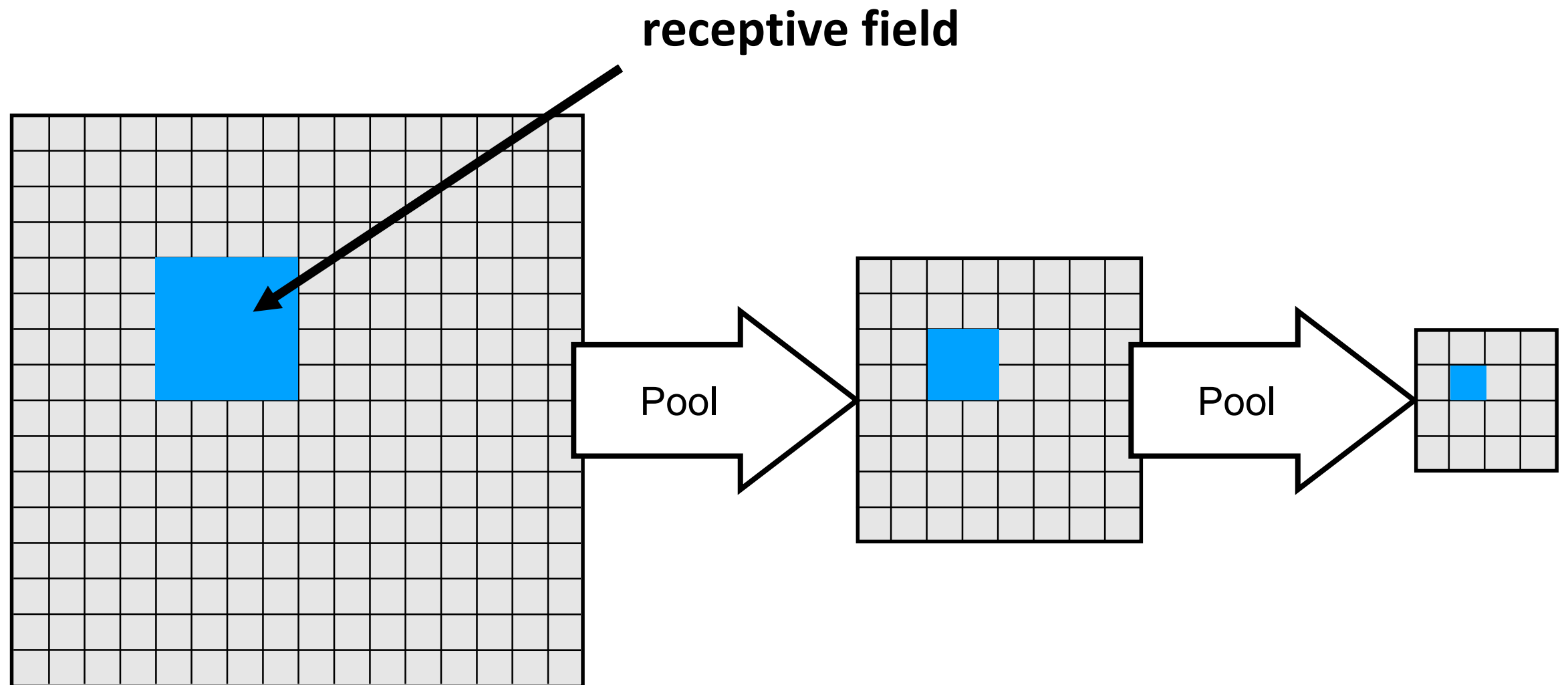
Further Benefit of (Max-)Pooling



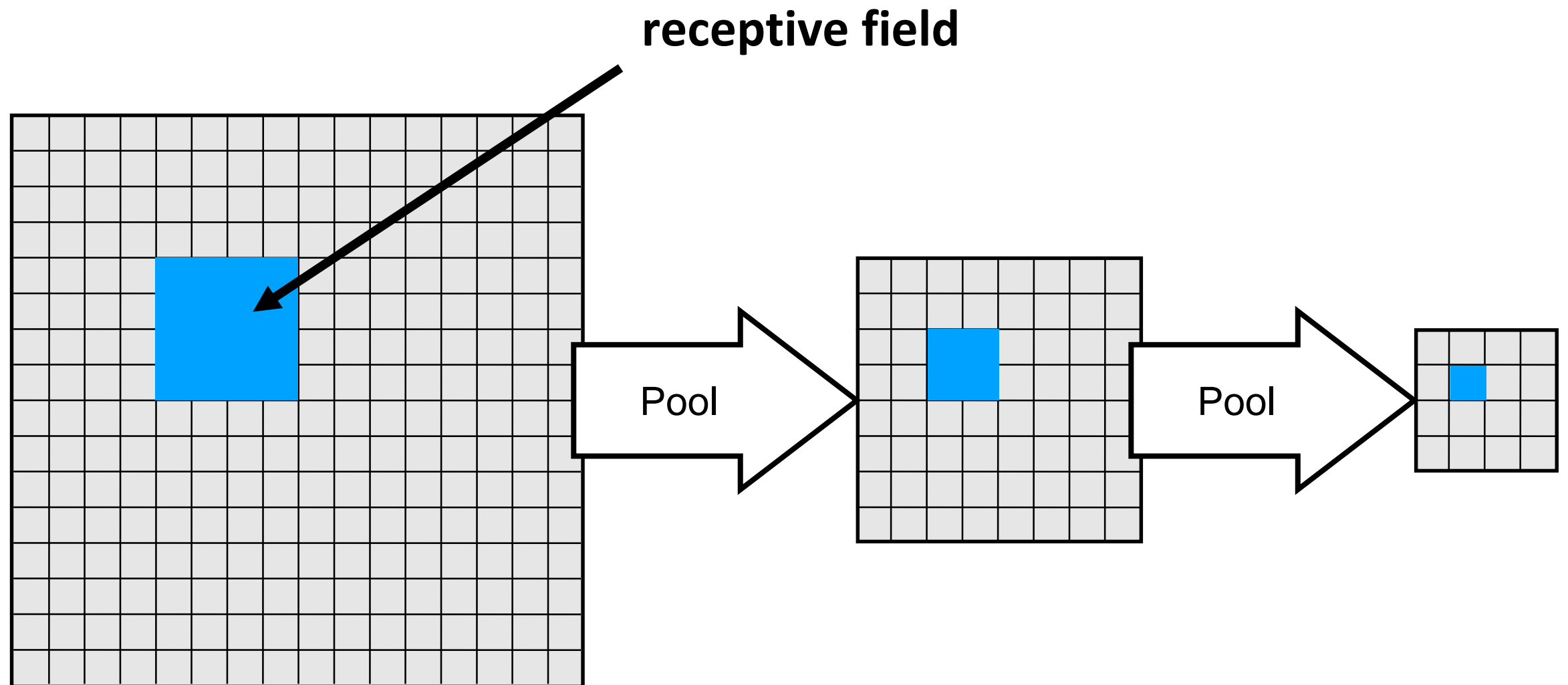
Further Benefit of (Max-)Pooling



Further Benefit of (Max-)Pooling

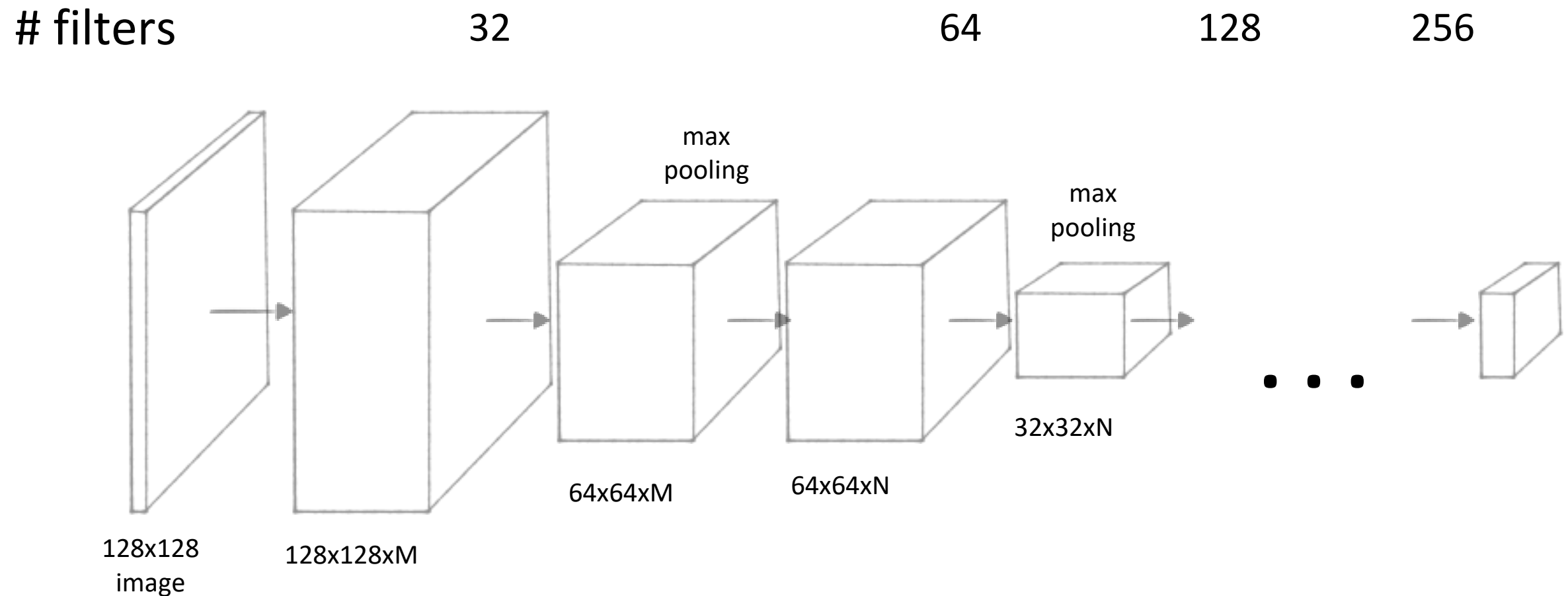


Further Benefit of (Max-)Pooling

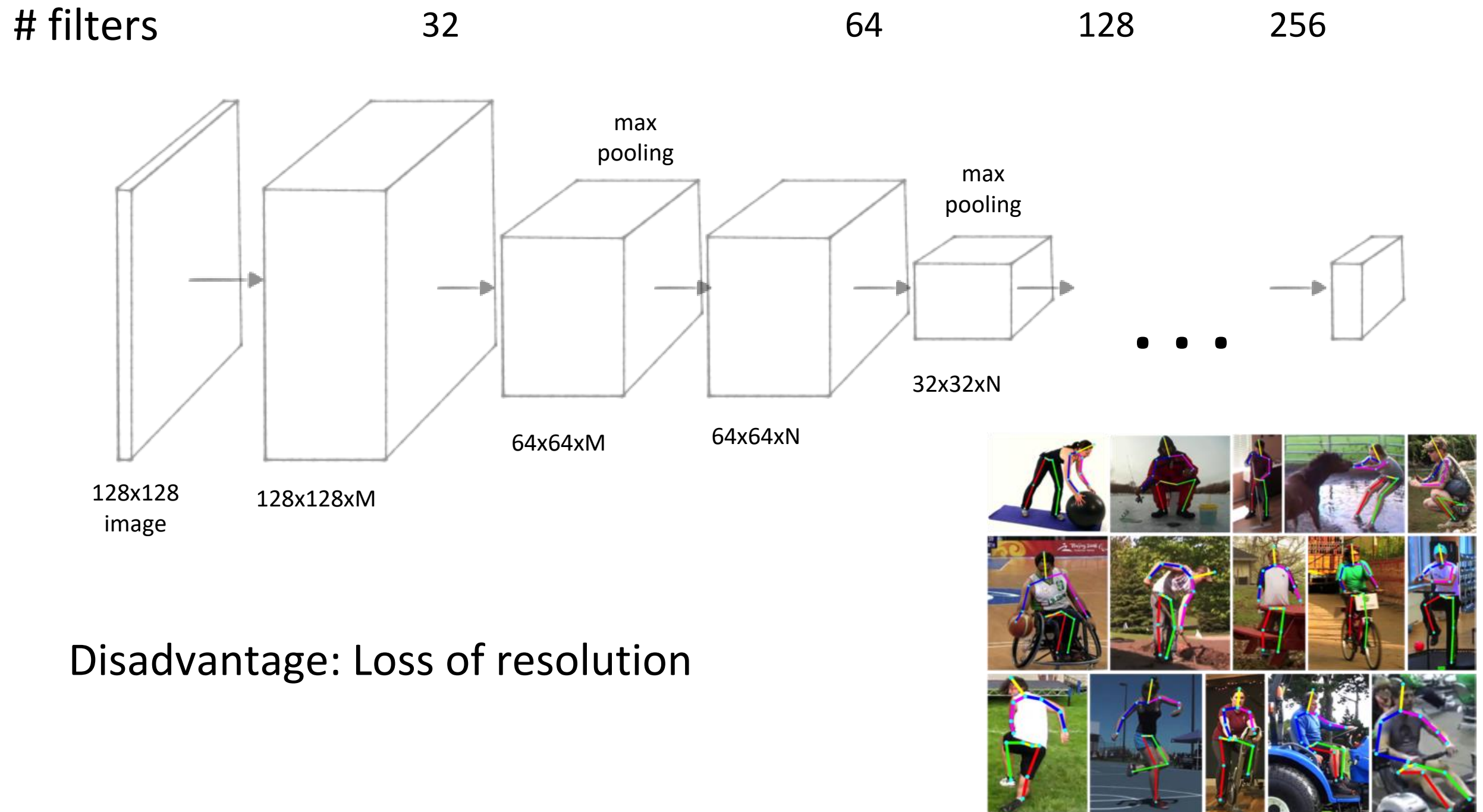


Pooling increases receptive field size!

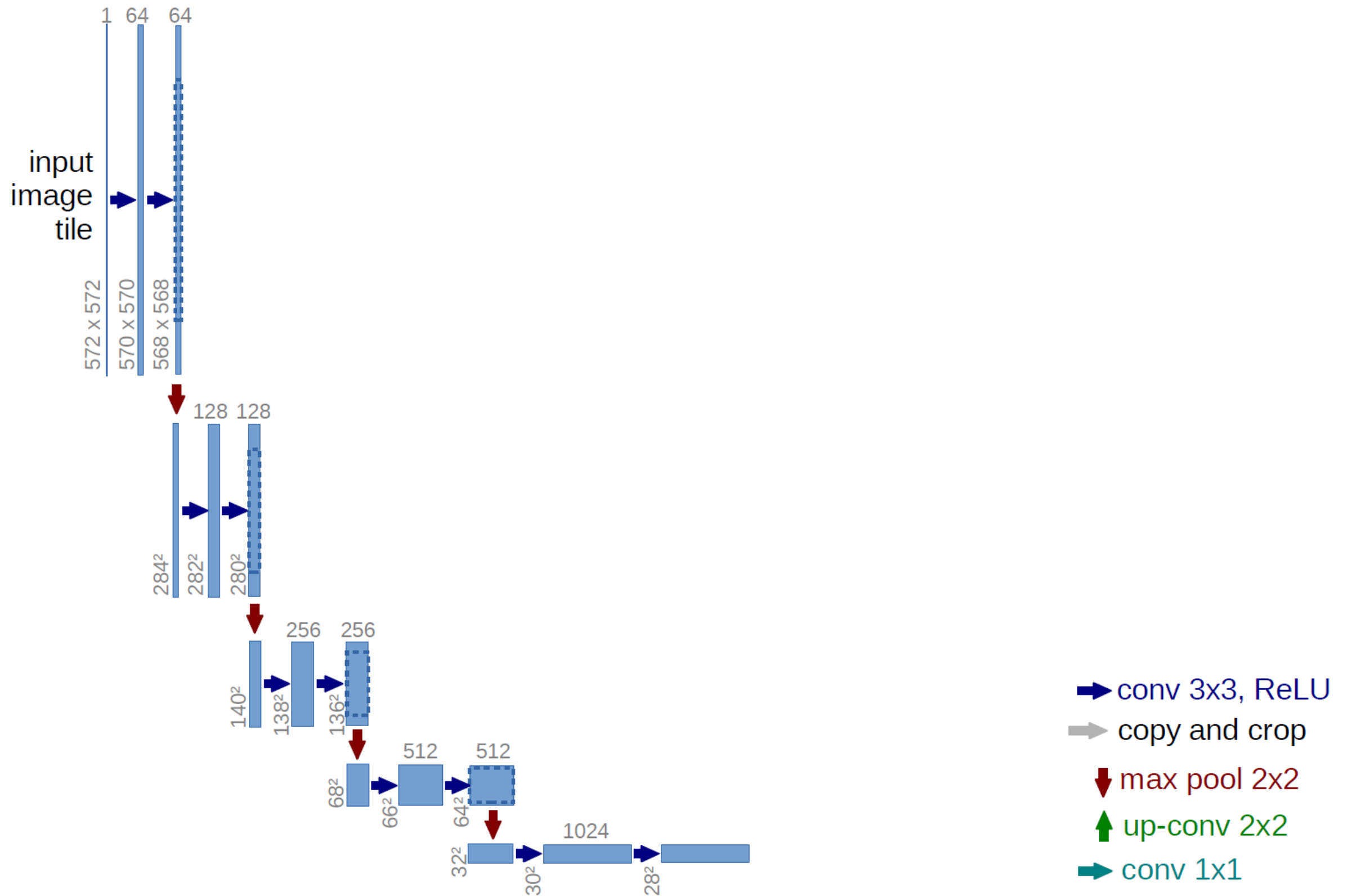
Fully-Convolutional Network for Detection



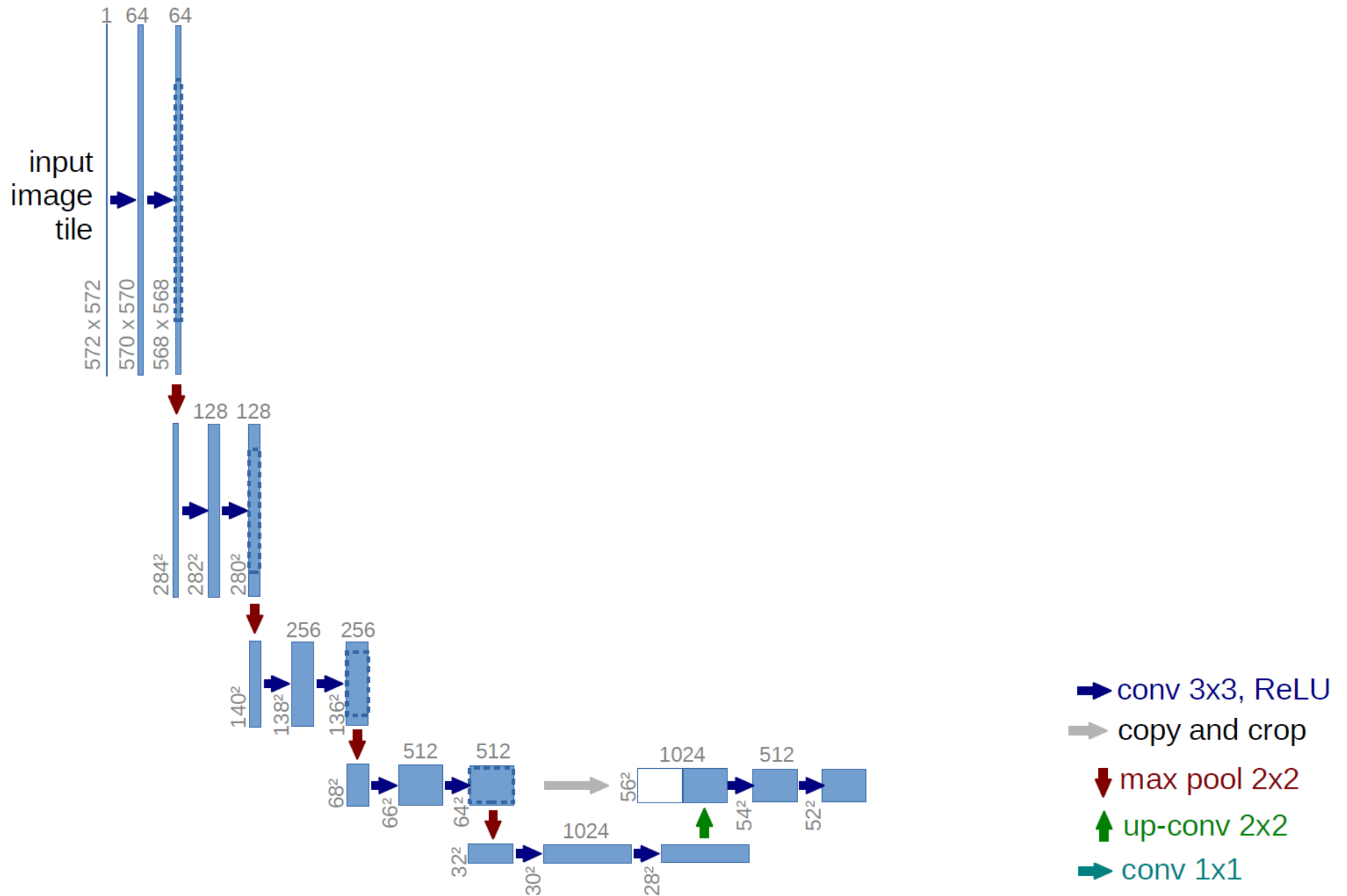
Fully-Convolutional Network for Detection



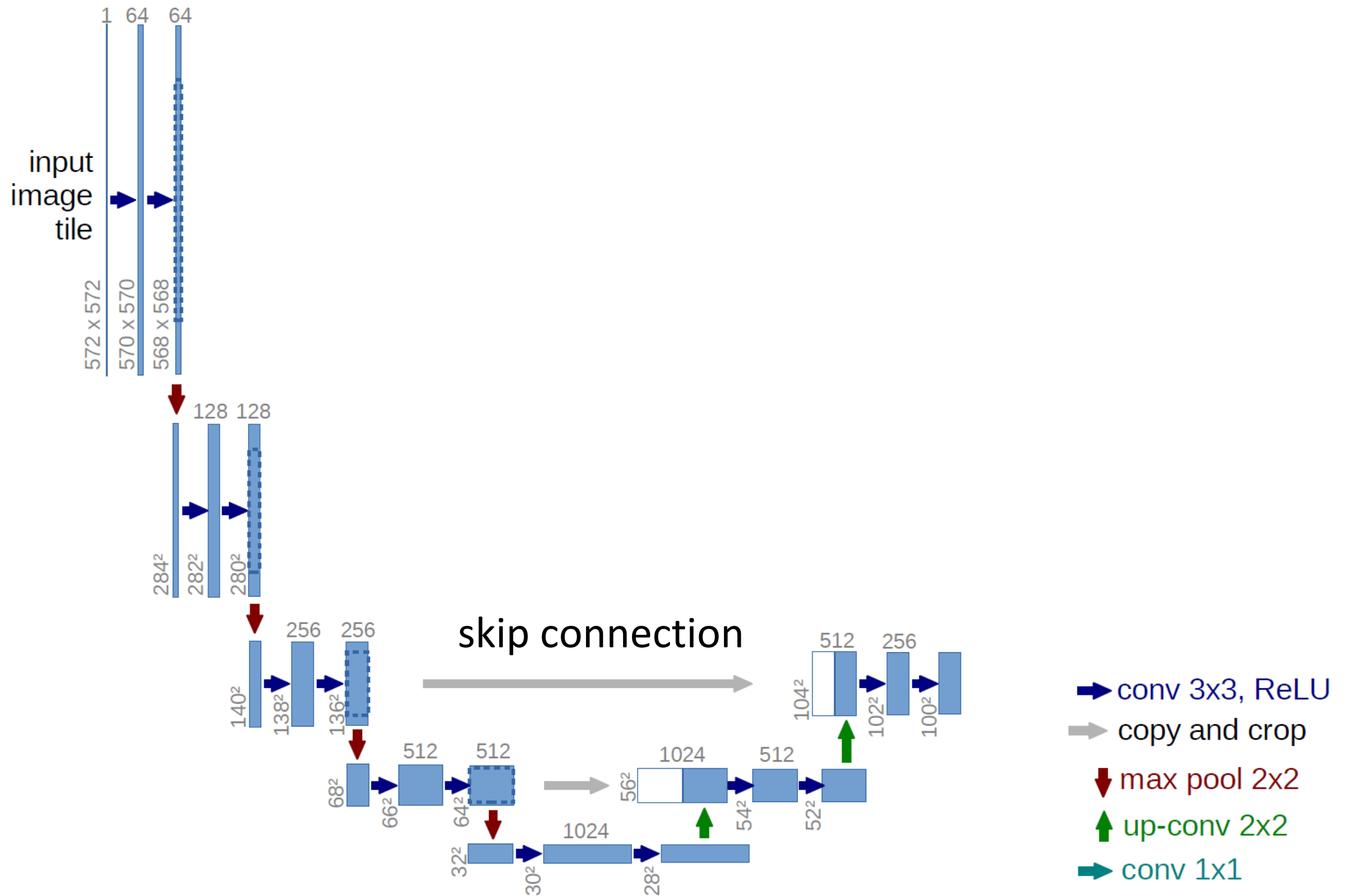
U-Net: Decrease and Increase Resolution



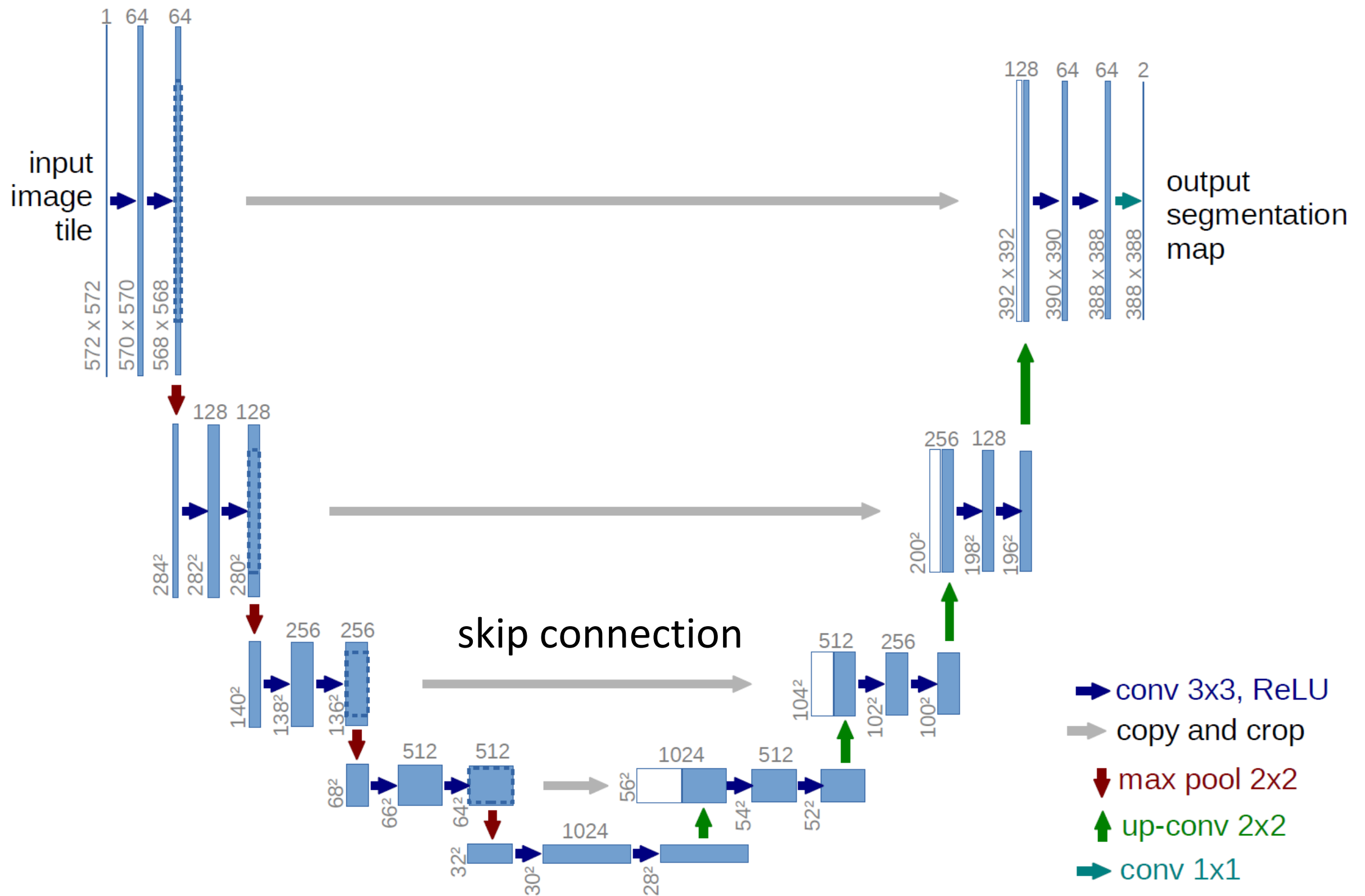
U-Net: Decrease and Increase Resolution



U-Net: Decrease and Increase Resolution

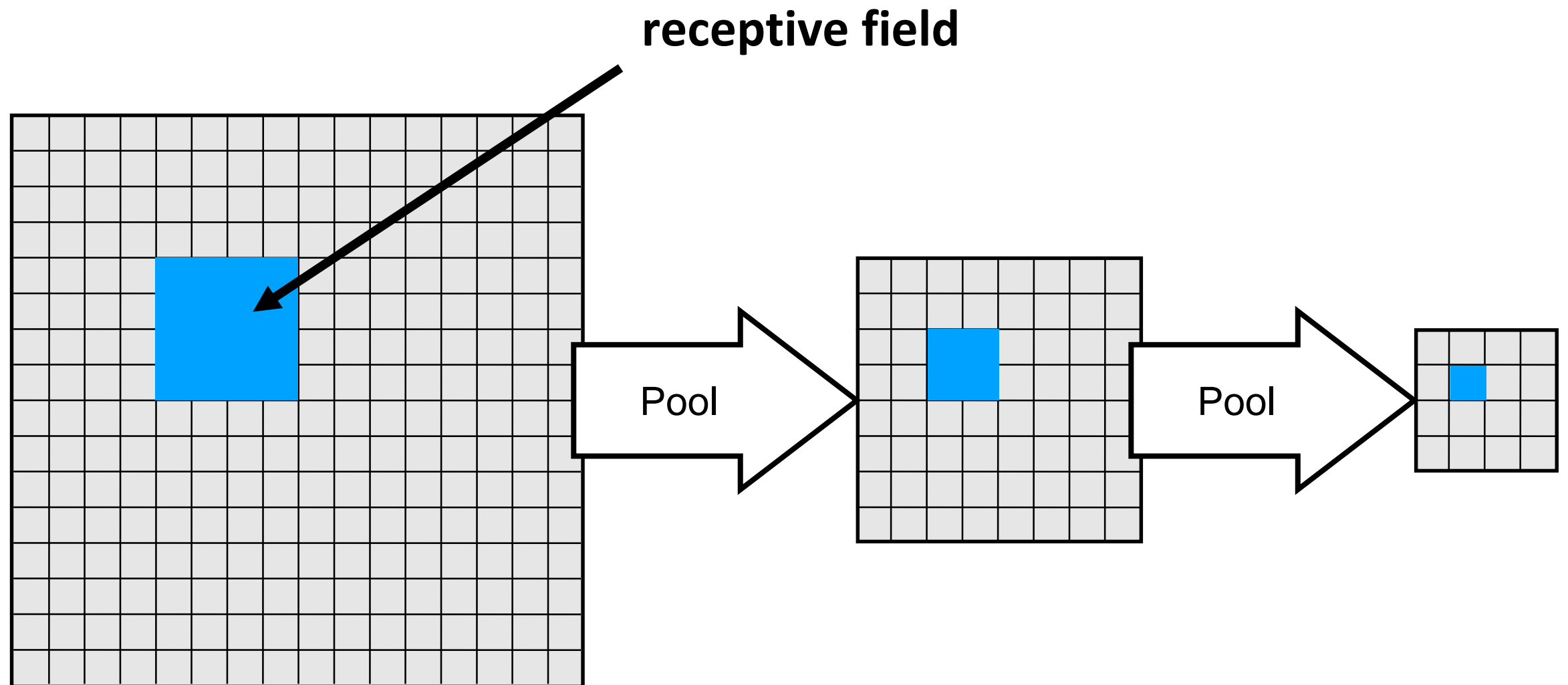


U-Net: Decrease and Increase Resolution



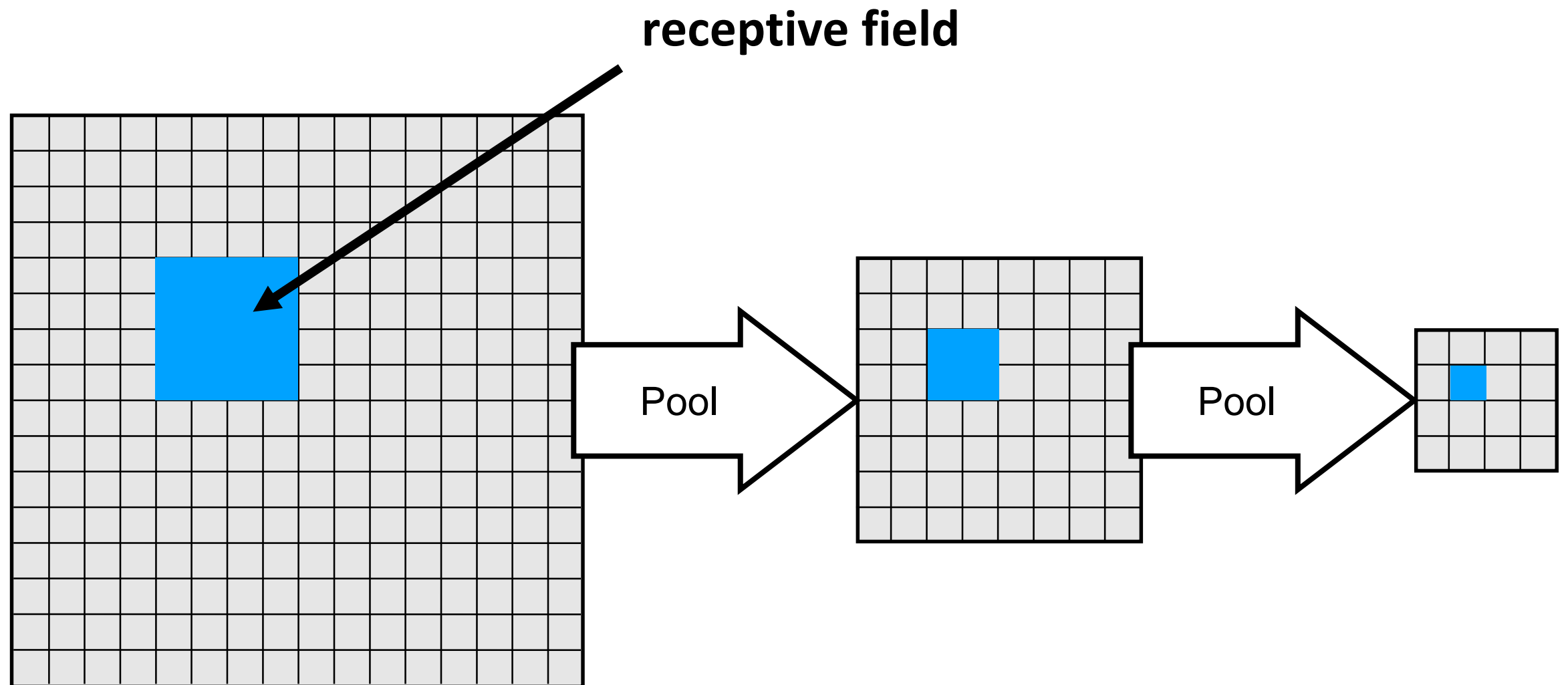
Convolutional Layers

Increasing the Receptive Field by (Max-)Pooling



Advantage: Introduces robustness to small shifts

Increasing the Receptive Field by (Max-)Pooling



Advantage: Introduces robustness to small shifts

... but we might not always want that

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

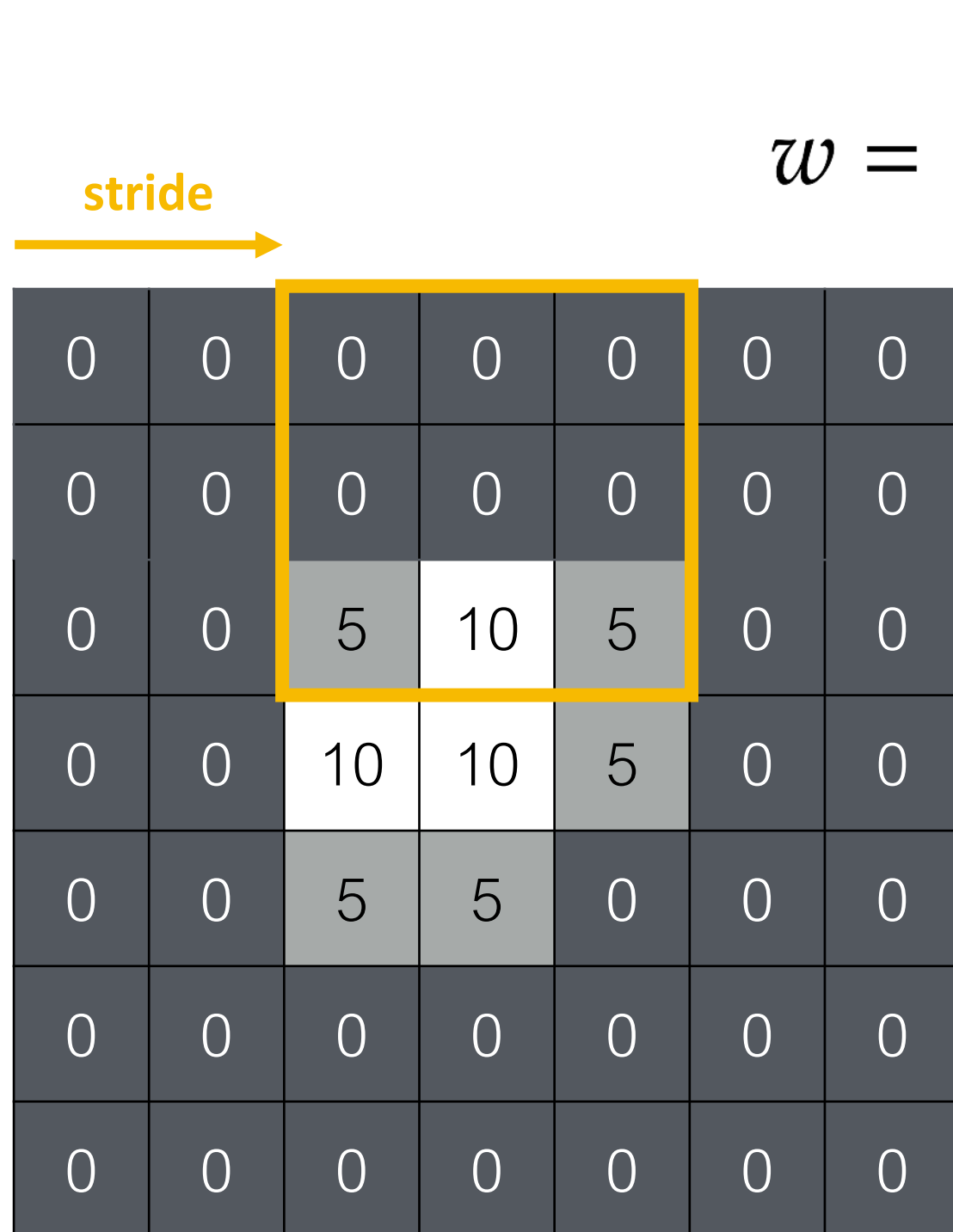
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0		

$I \star w$

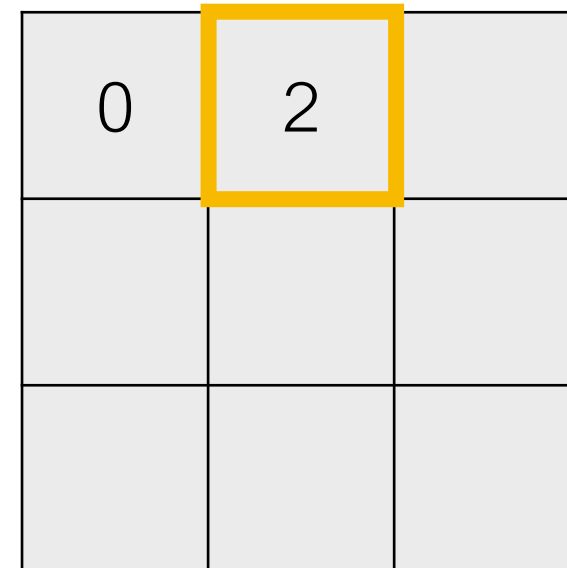
Learnable Alternative: Strided Filtering



I

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0



$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0	2	0

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0	2	0
2		

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0	2	0
2	8	

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0	2	0
2	8	1

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

0	2	0
2	8	1
0		

$I \star w$

Learnable Alternative: Strided Filtering

$$w = \frac{1}{5}$$

0	1	0
1	1	1
0	1	0

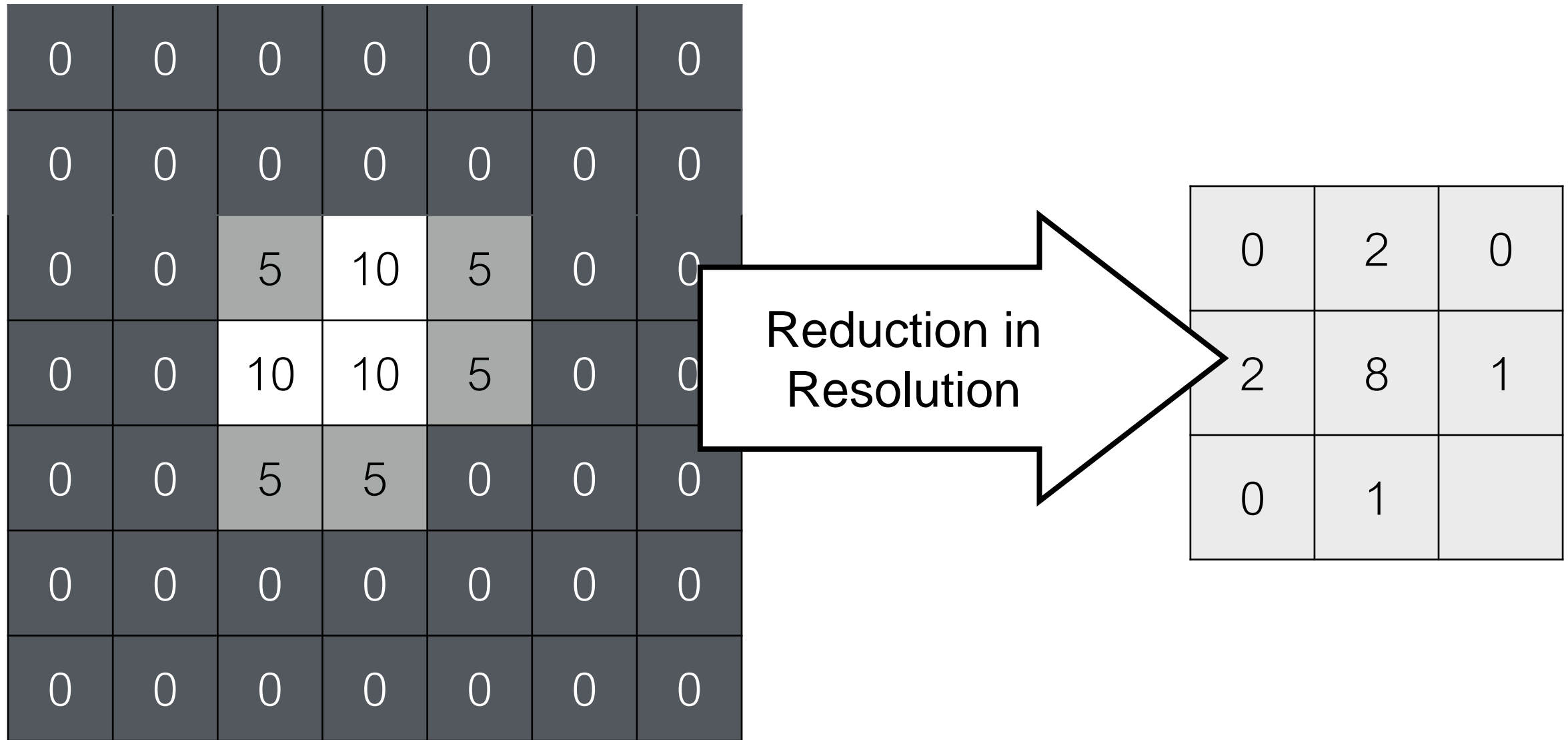
0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	5	10	5	0	0
0	0	10	10	5	0	0
0	0	5	5	0	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

I

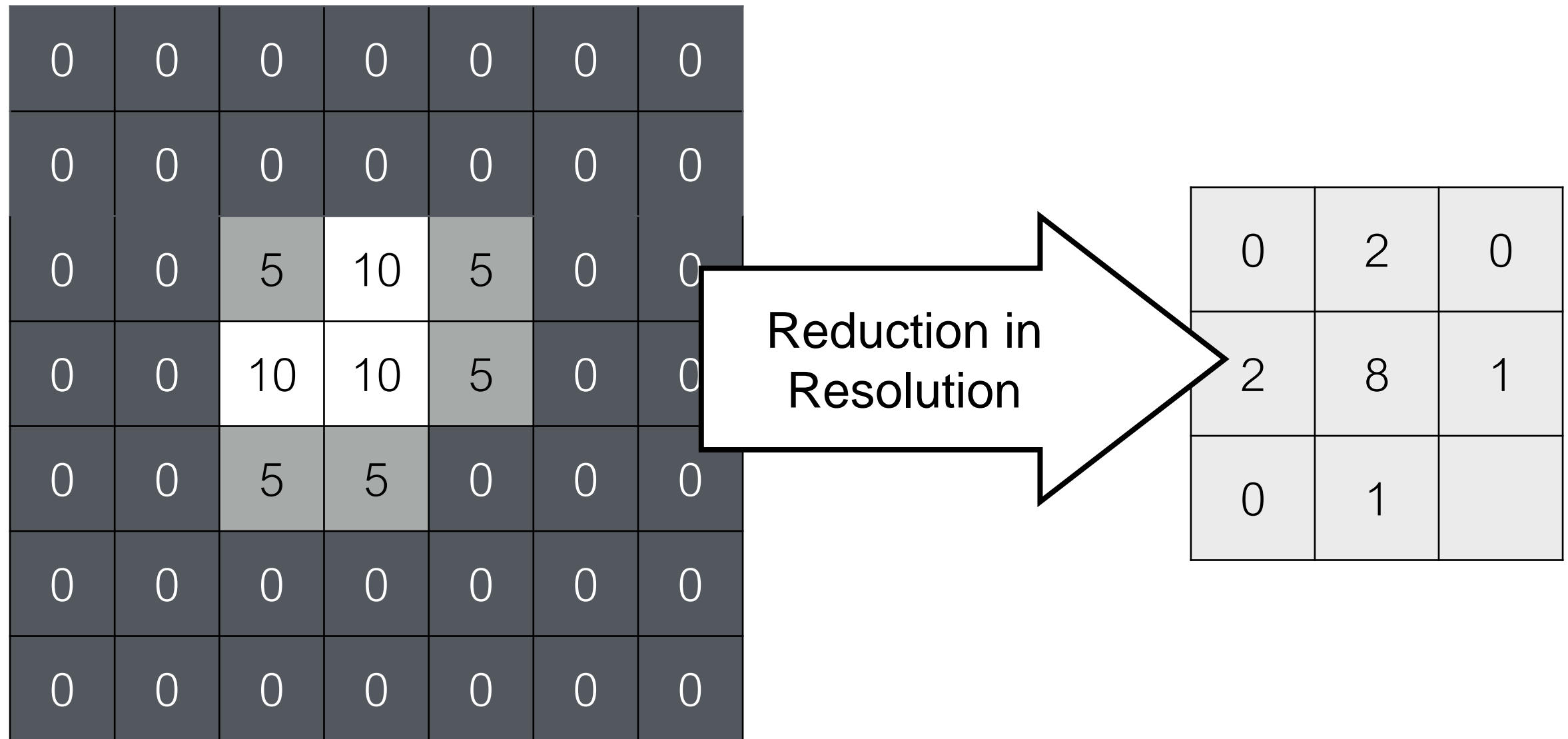
0	2	0
2	8	1
0	1	

$I \star w$

Strided Convolutions

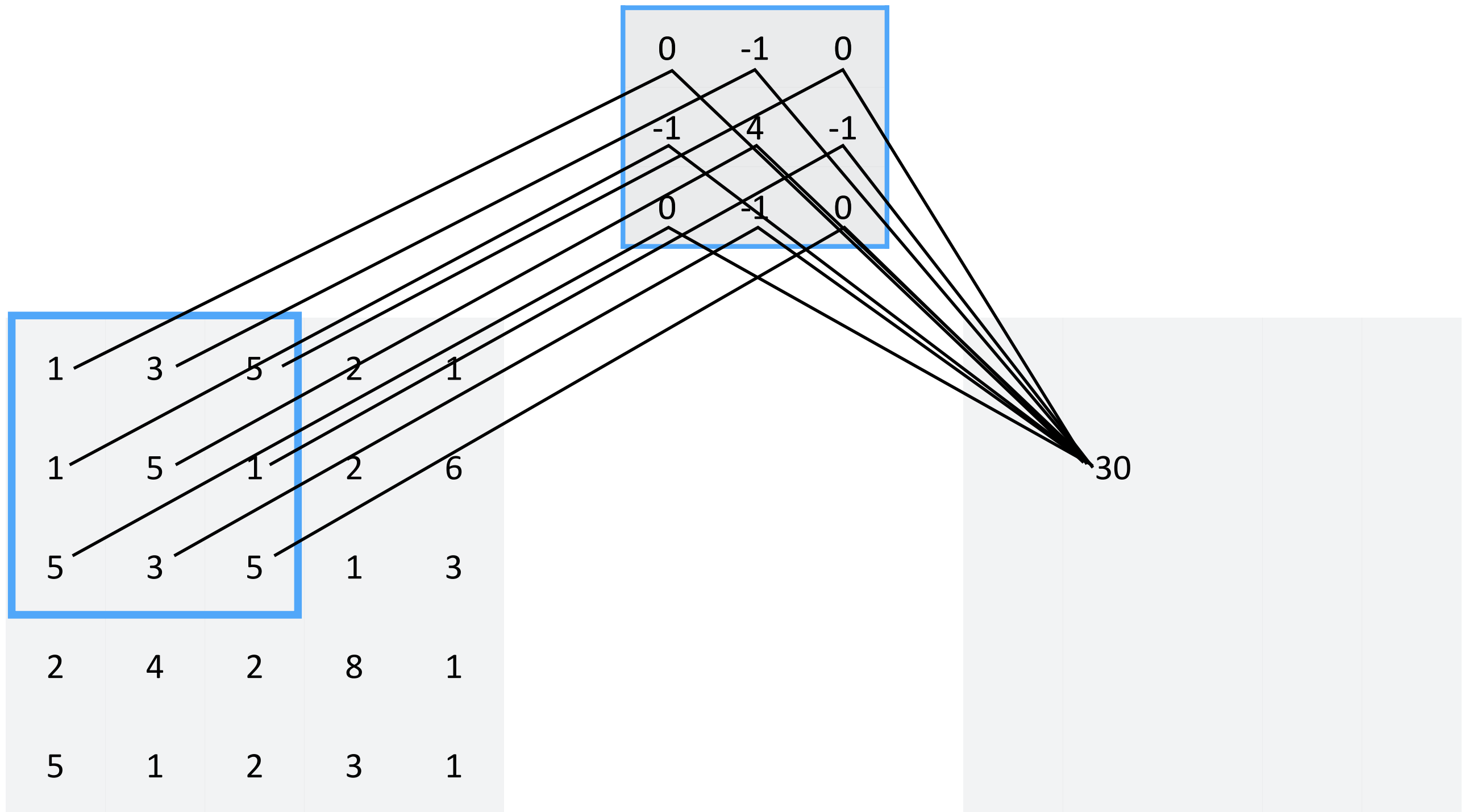


Strided Convolutions



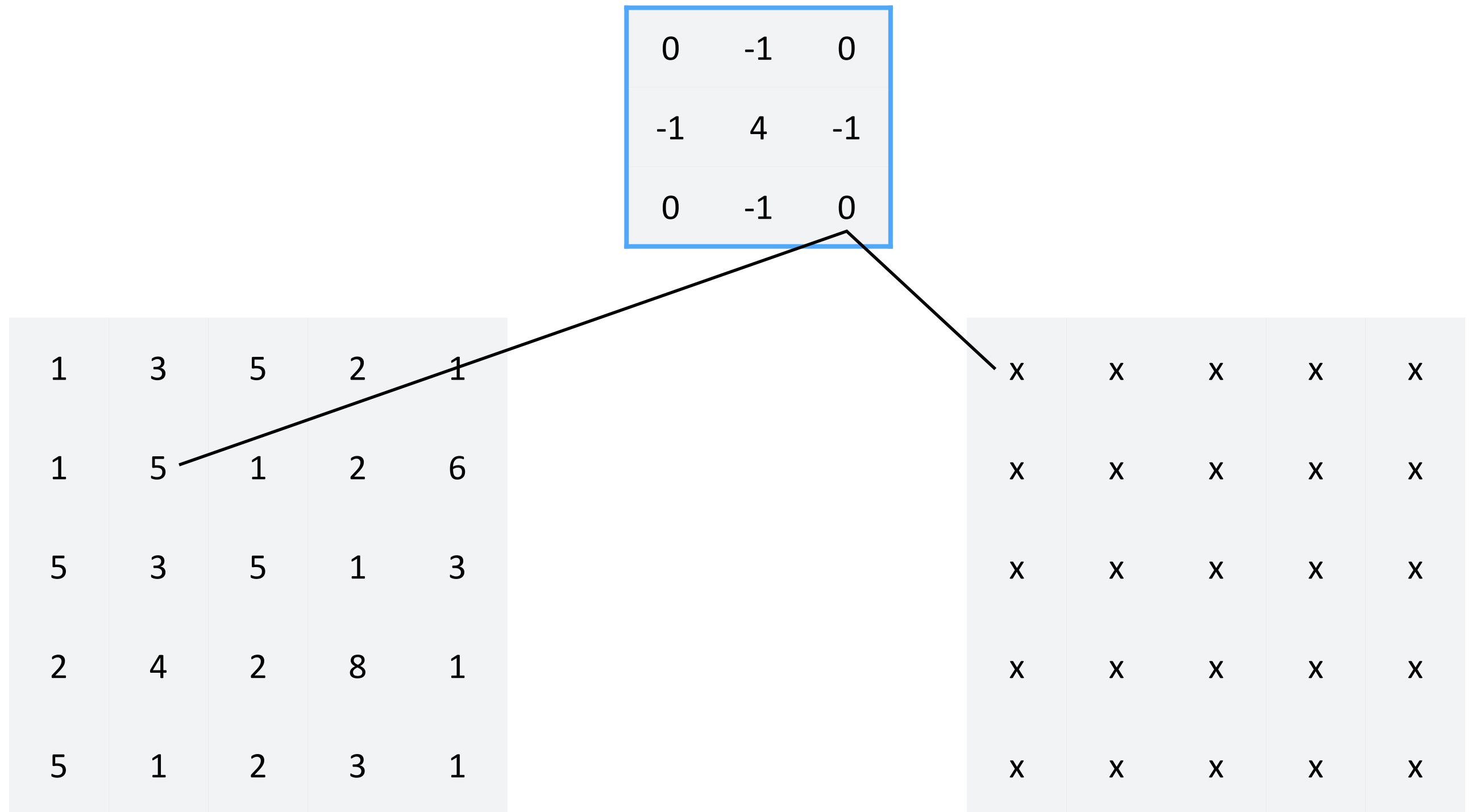
What if we want to **increase** resolution?

Different View on Strided Convolutions



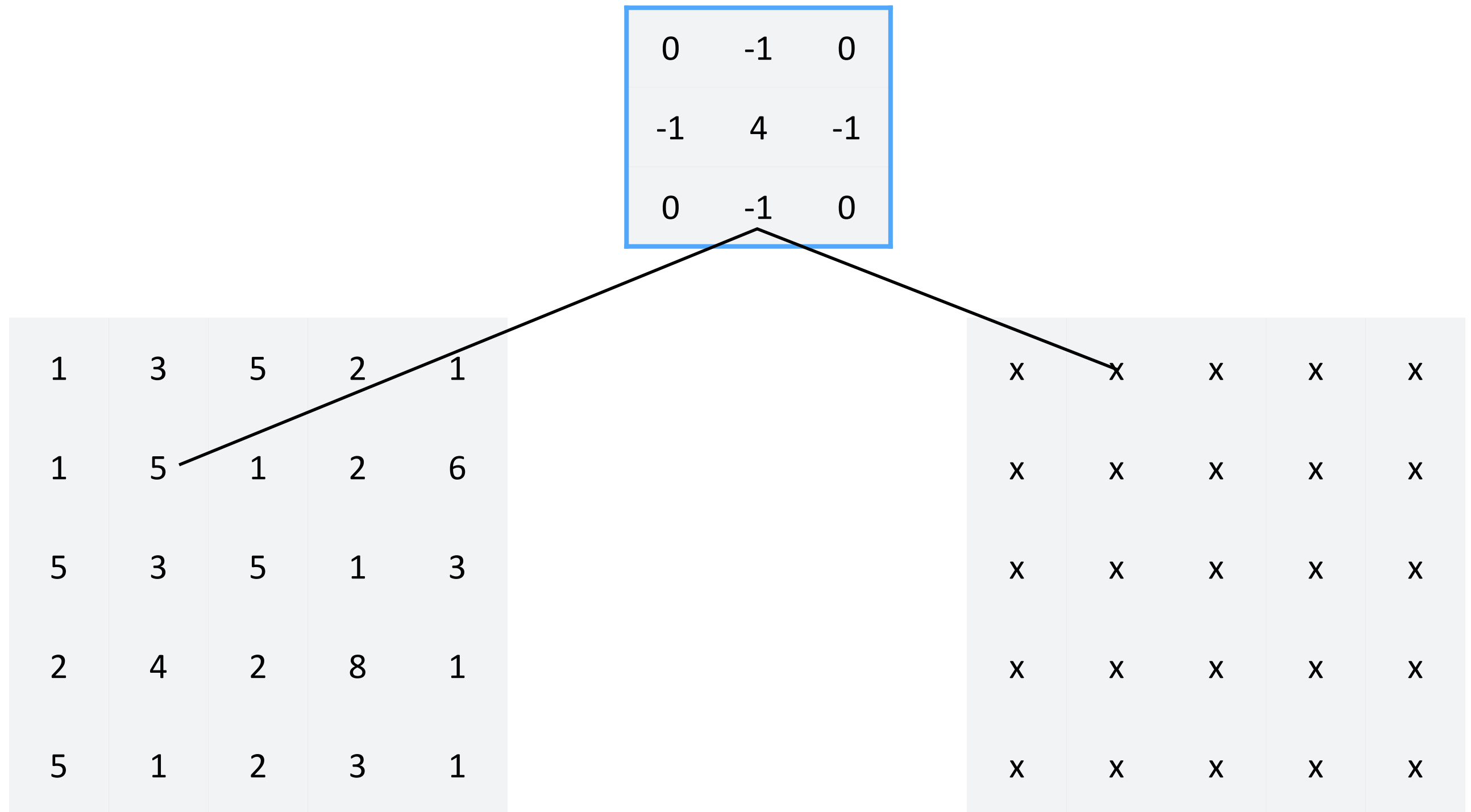
Standard view: Multiple pixels contribute to one pixel

Different View on Strided Convolutions



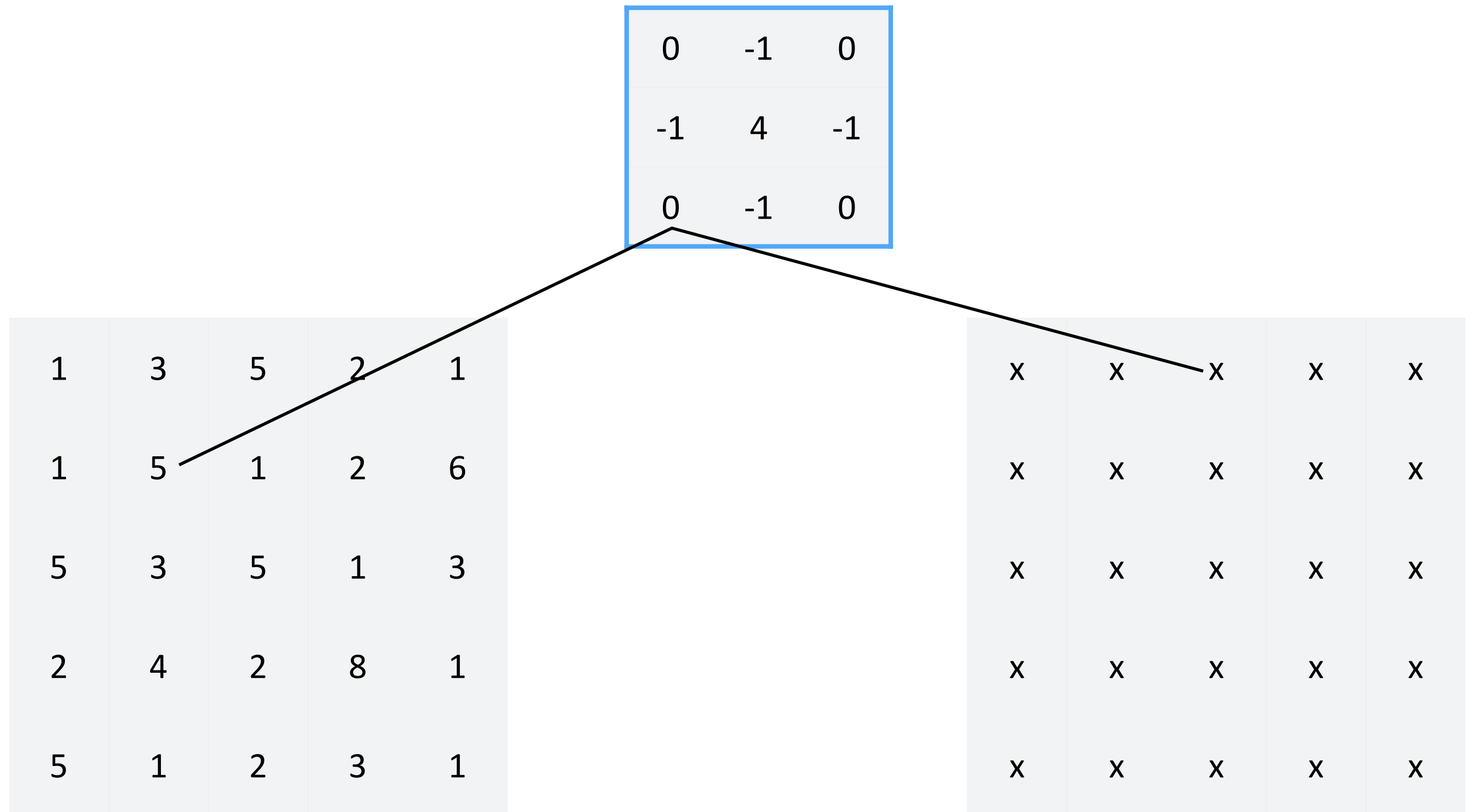
Or: One pixel contributes to many pixels

Different View on Strided Convolutions



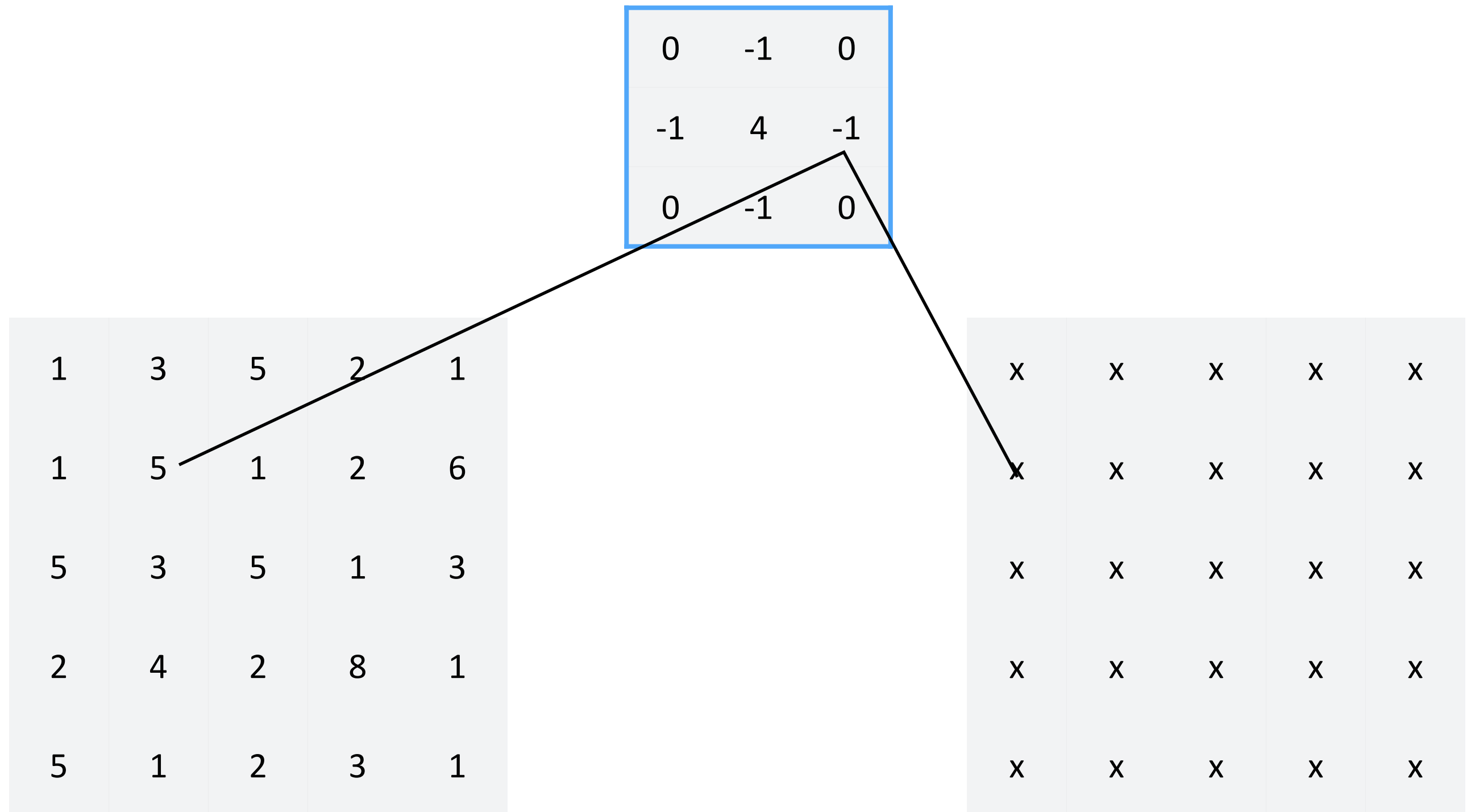
Or: One pixel contributes to many pixels

Different View on Strided Convolutions



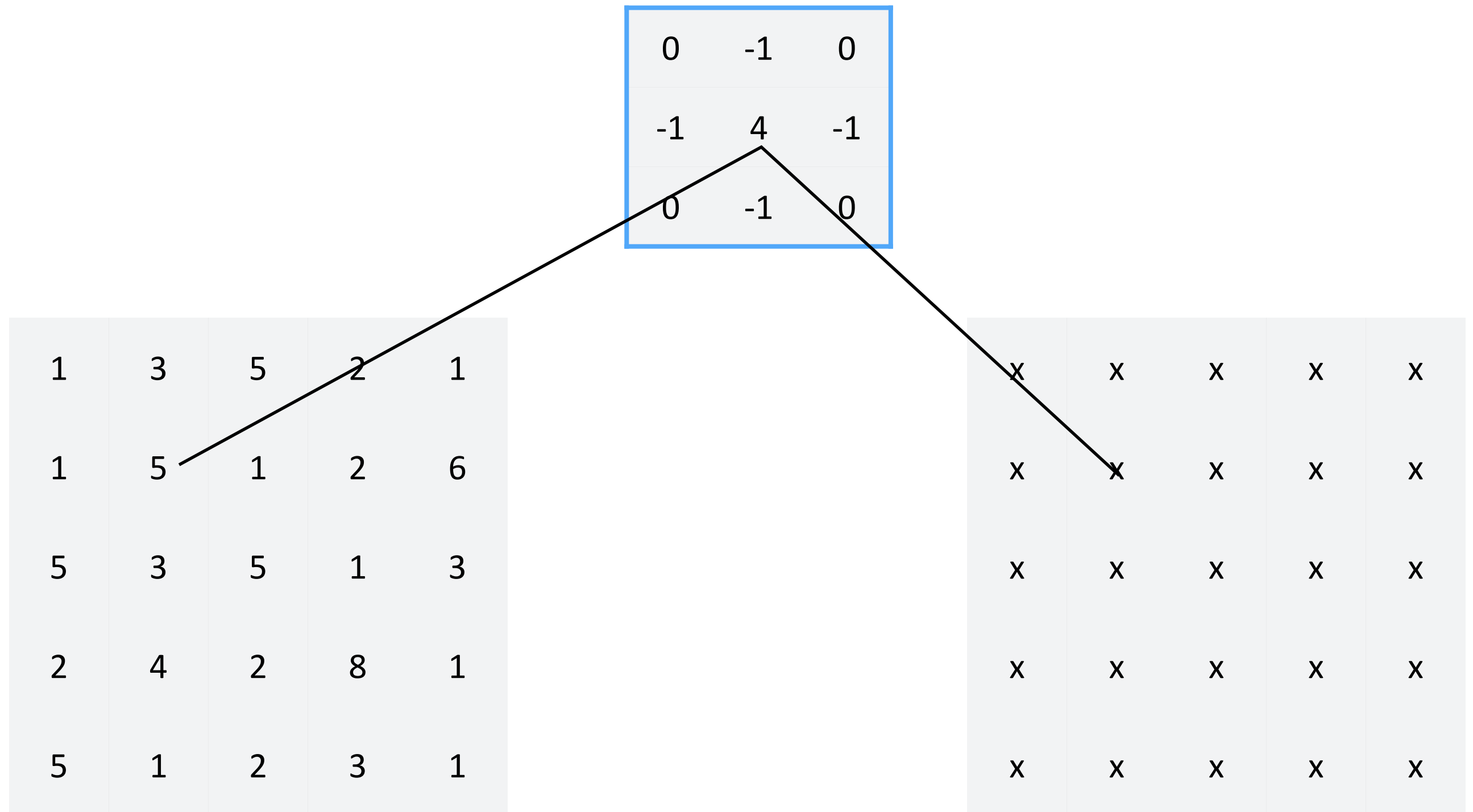
Or: One pixel contributes to many pixels

Different View on Strided Convolutions



Or: One pixel contributes to many pixels

Different View on Strided Convolutions



Or: One pixel contributes to many pixels

Different View on Strided Convolutions

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

x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

Or: One pixel contributes to many pixels

Different View on Strided Convolutions

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

x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

**Transposed
convolution**

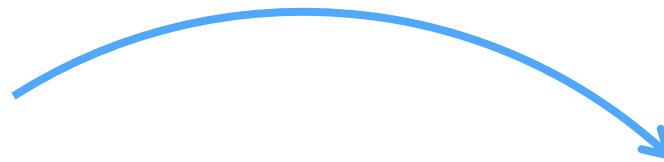
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

Or: One pixel contributes to many pixels

What's the point?

x	x	x
x	x	x
x	x	x

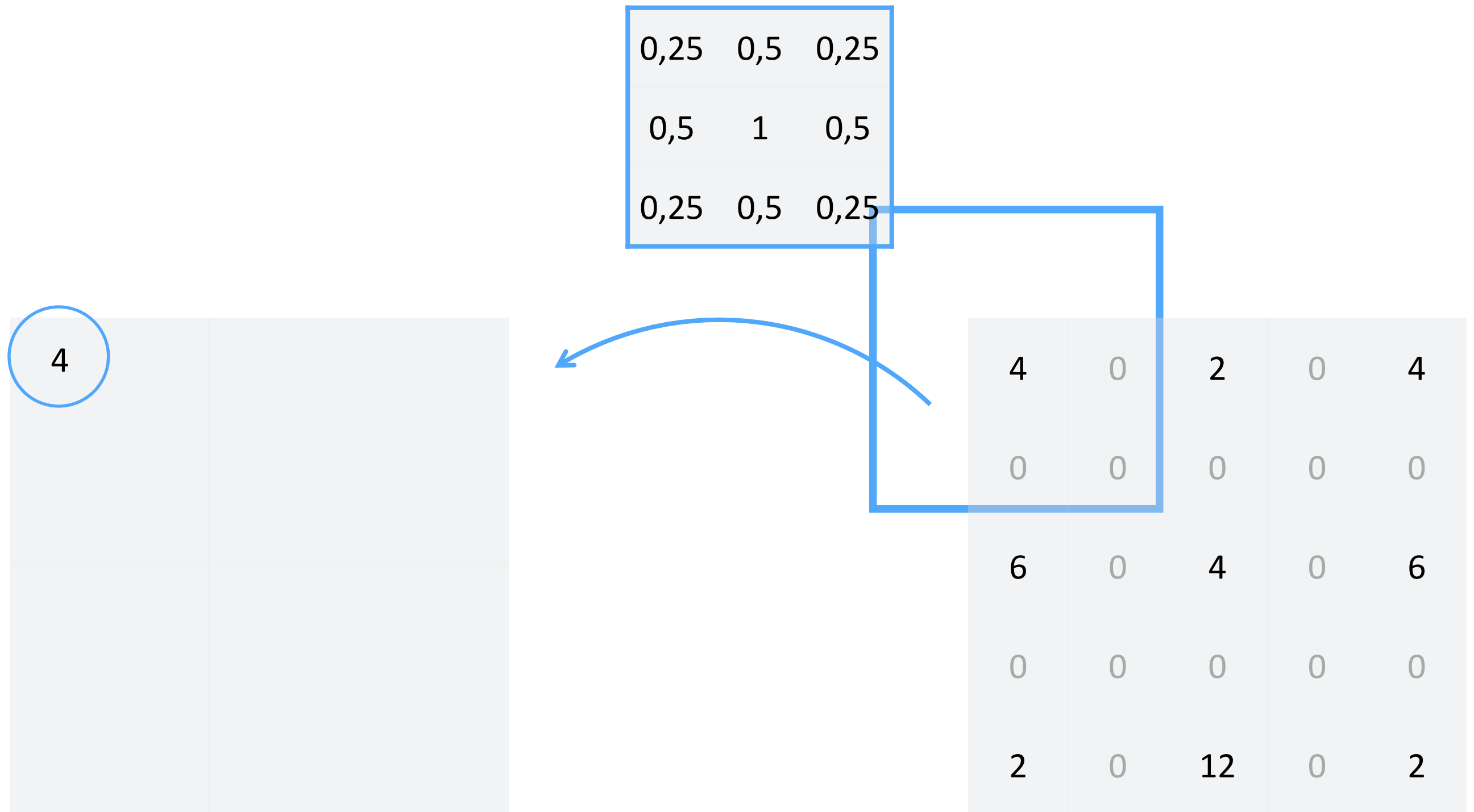
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x
x	x	x	x	x



4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Strided filtering reduces resolution

What's the point?



Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3		

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	3

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	3	4

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	3	4
5				

	4	0	2	0	4
	0	0	0	0	0
	6	0	4	0	6
	0	0	0	0	0
	2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	3	4
5	4			

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

Transposed strided filtering increases resolution

What's the point?

0,25	0,5	0,25
0,5	1	0,5
0,25	0,5	0,25

4	3	2	3	4
5	4	3	4	5
6	5	4	5	6
4	6	8	6	4
2	7	12	7	2

4	0	2	0	4
0	0	0	0	0
6	0	4	0	6
0	0	0	0	0
2	0	12	0	2

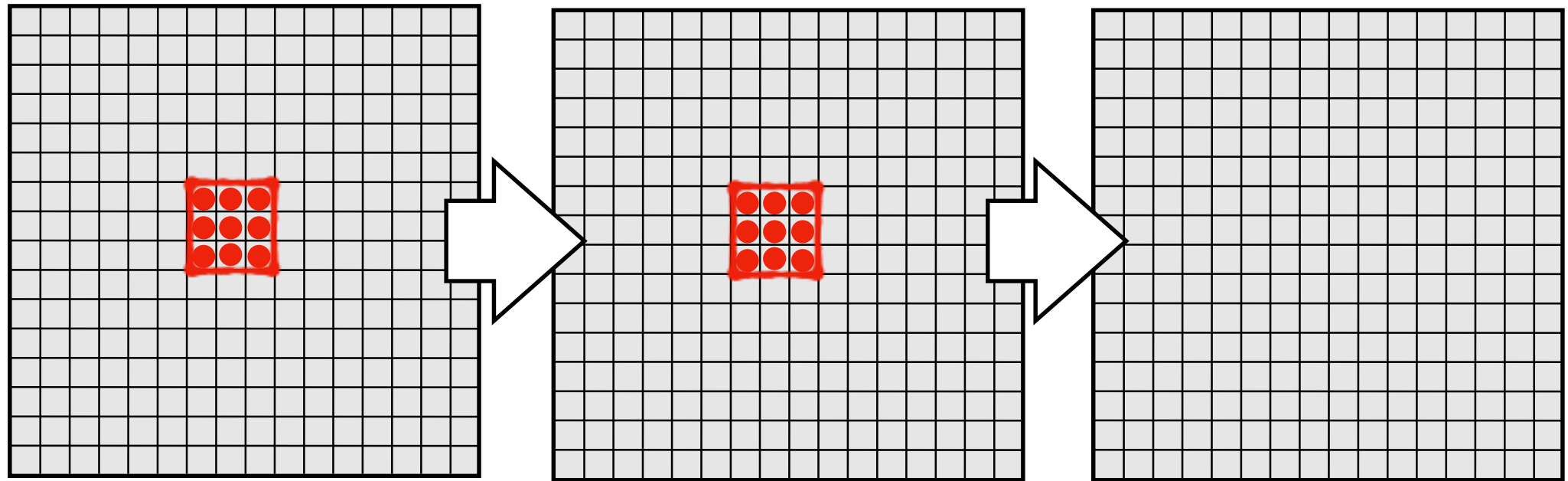
Transposed strided filtering increases resolution

Transposed Convolutions

- Not to be confused with matrix transpose
- Parameters are learned, not computed from the downsampling filtering stage
- Also known as:
 - fractionally strided convolutions
 - deconvolution (misleading, since transposed convolutions do not undo effect of convolution)

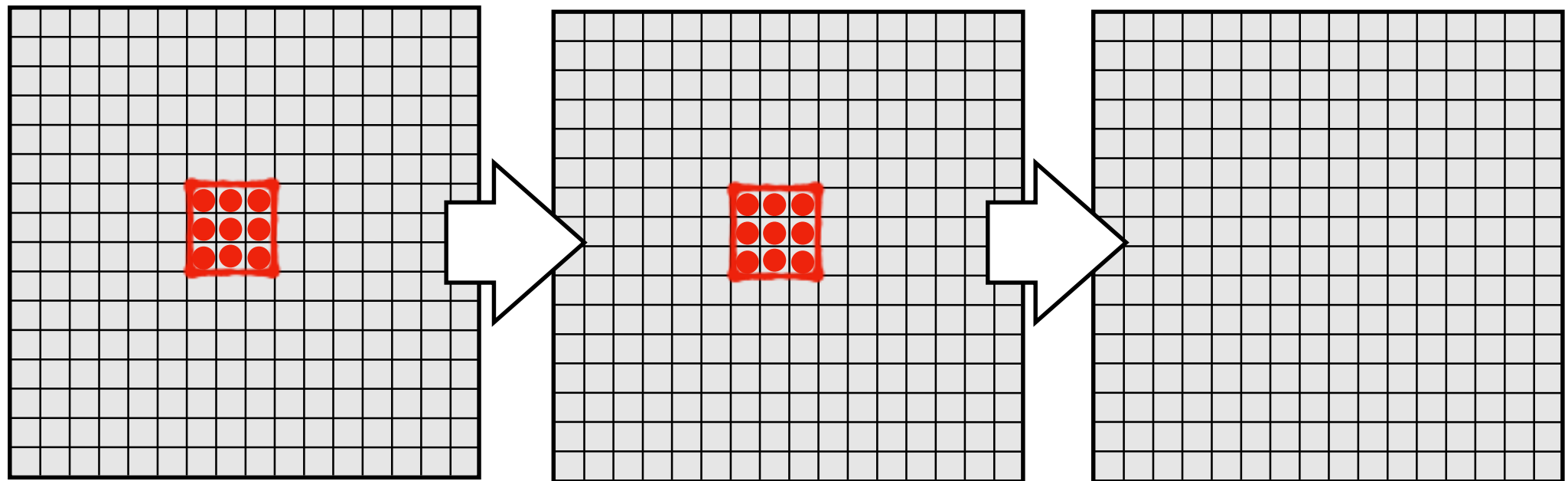
Receptive Field of Convolutions

Standard convolution:

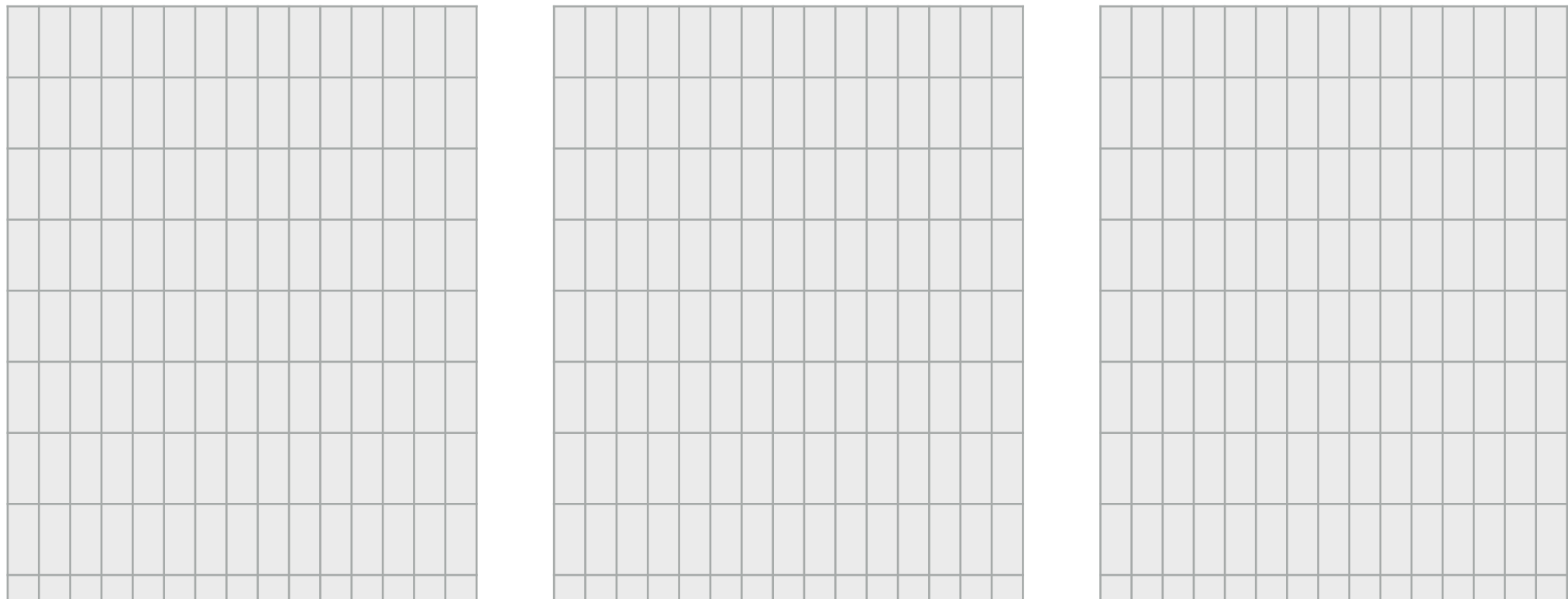


Receptive Field of Convolutions

Standard convolution:

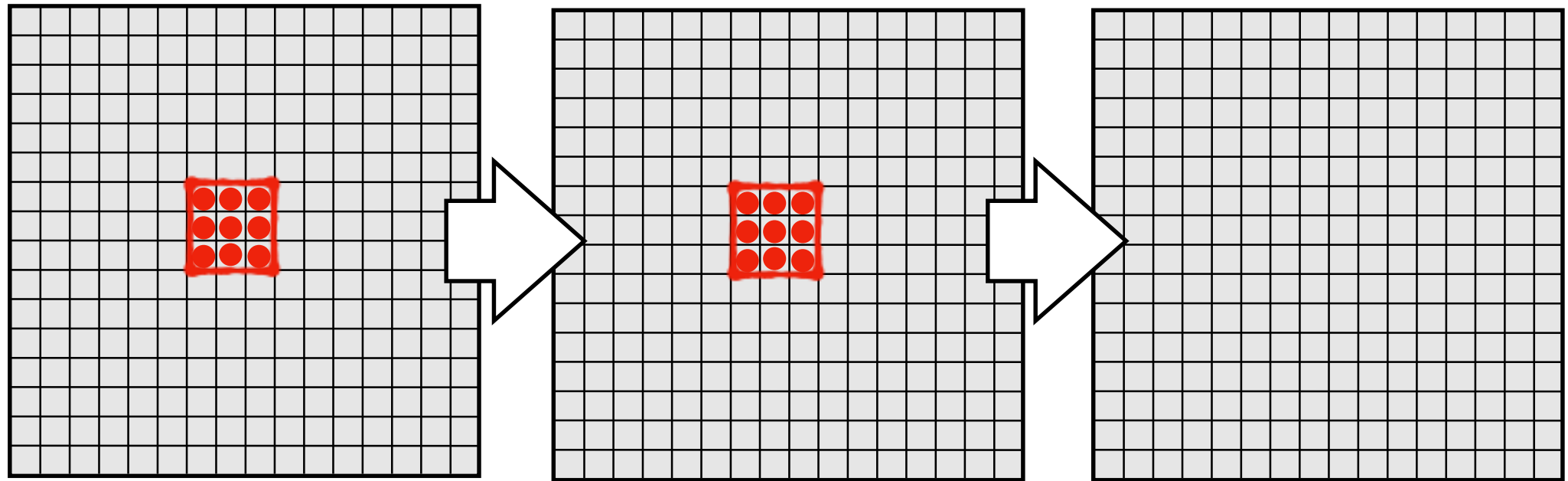


Receptive field:

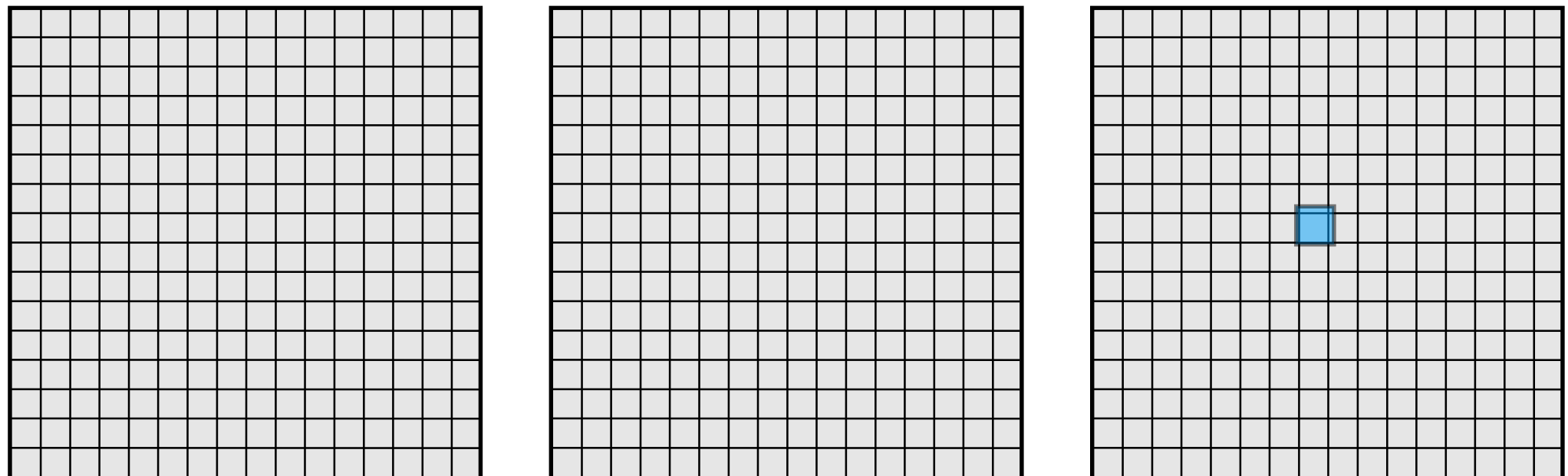


Receptive Field of Convolutions

Standard convolution:

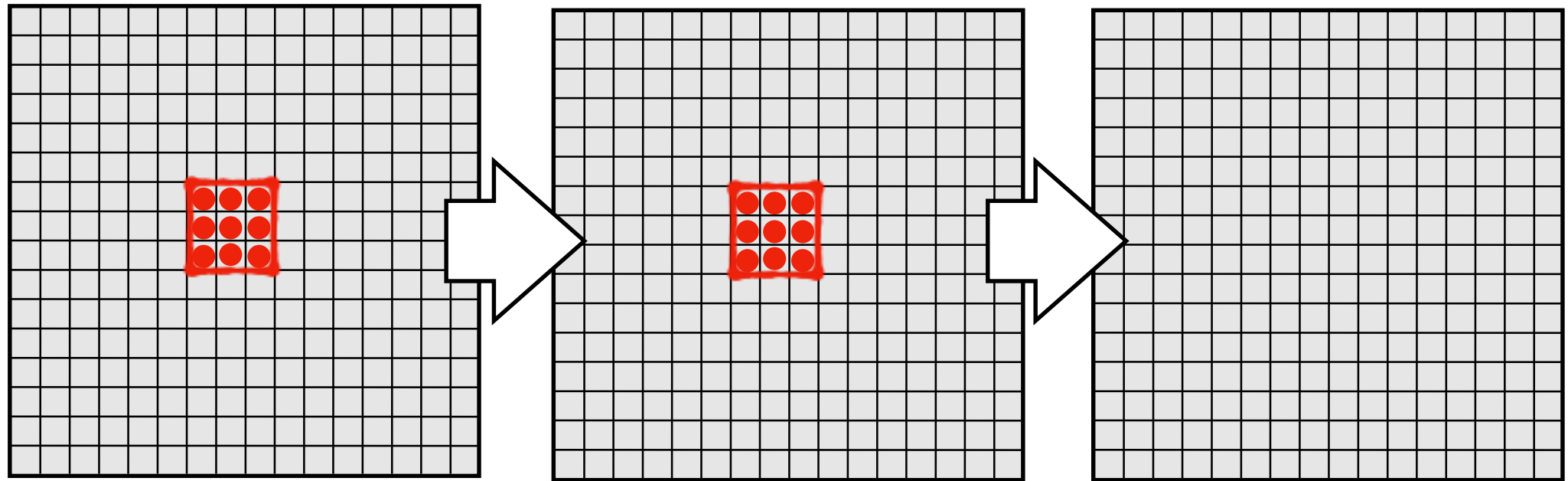


Receptive field:

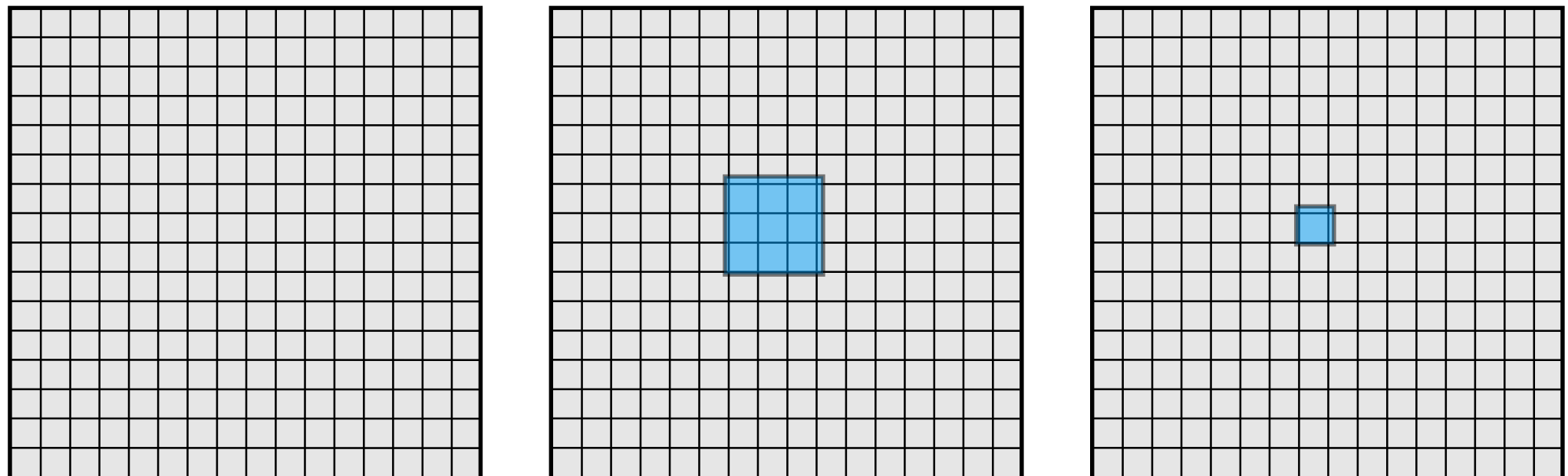


Receptive Field of Convolutions

Standard convolution:

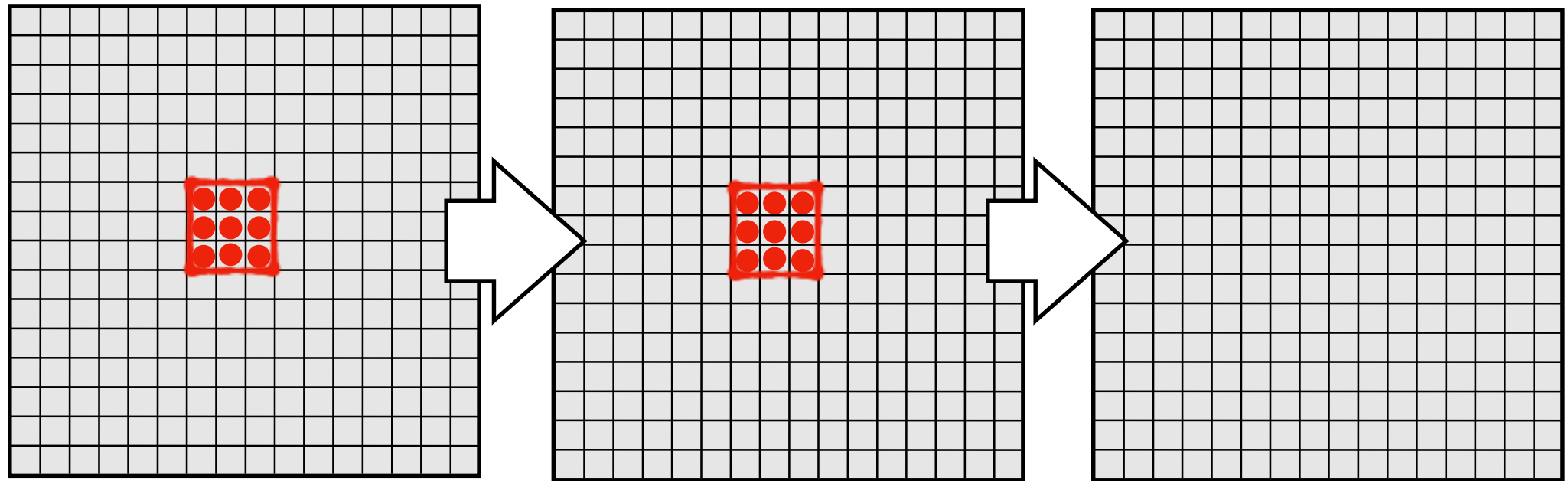


Receptive field:

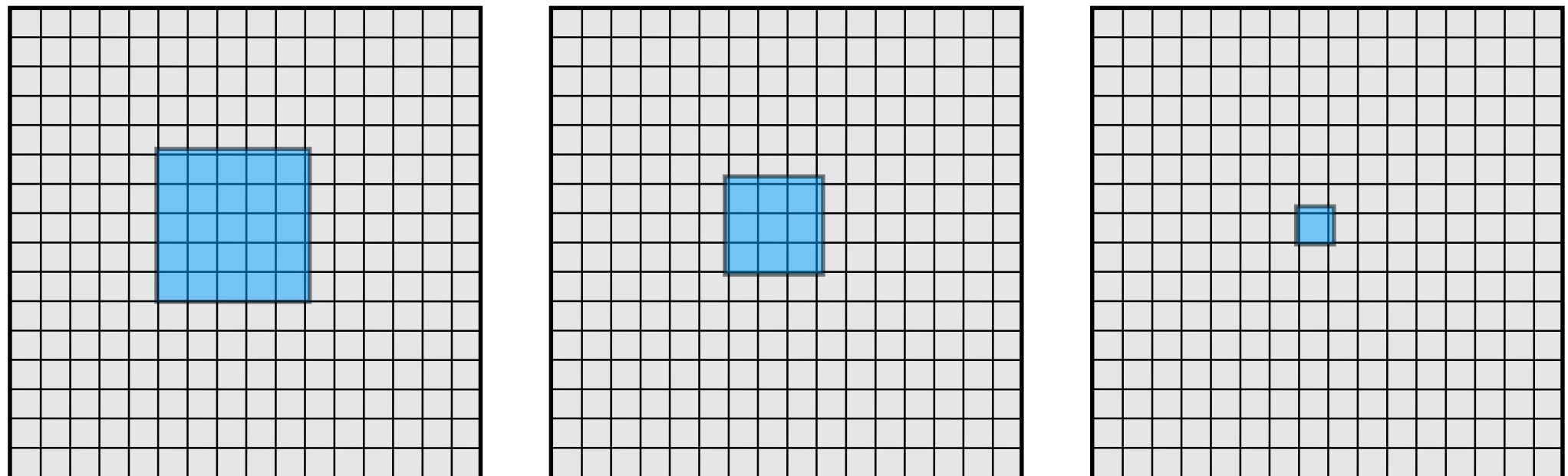


Receptive Field of Convolutions

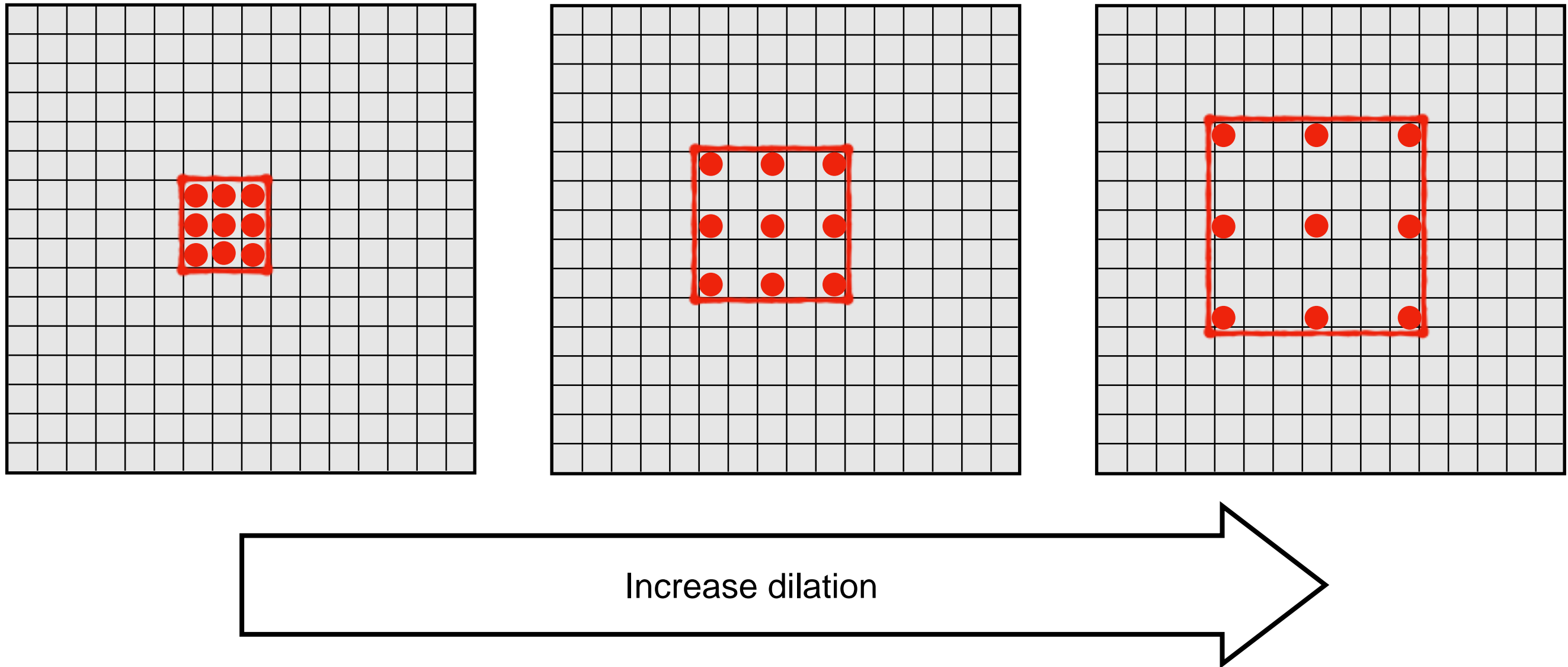
Standard convolution:



Receptive field:

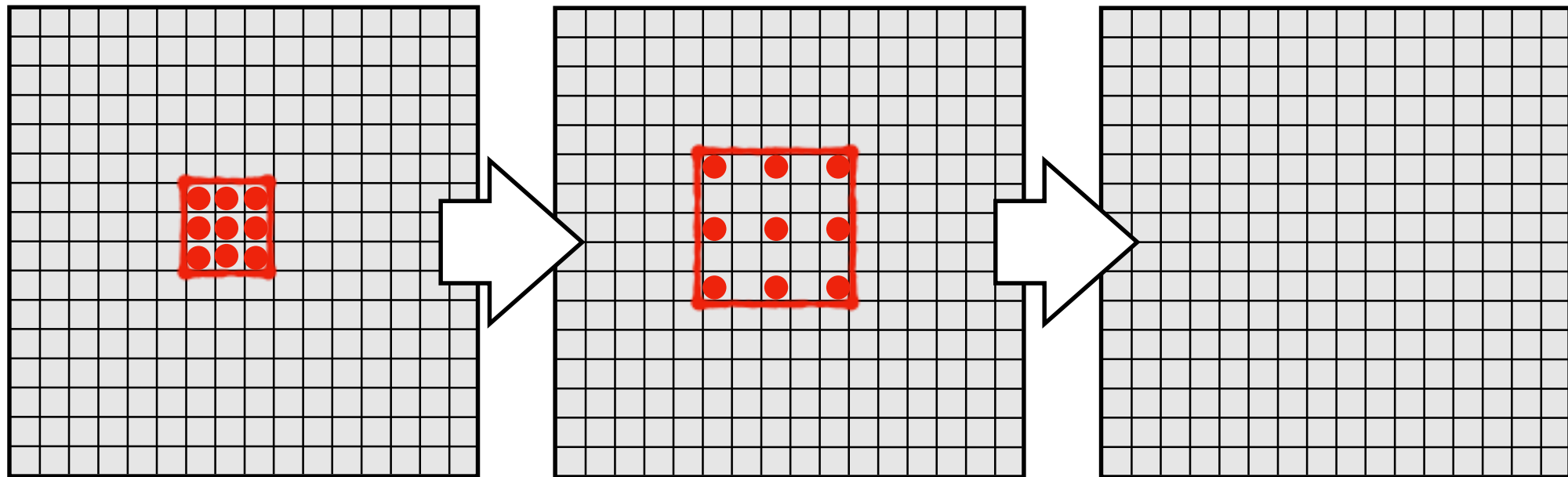


Dilated Convolutions



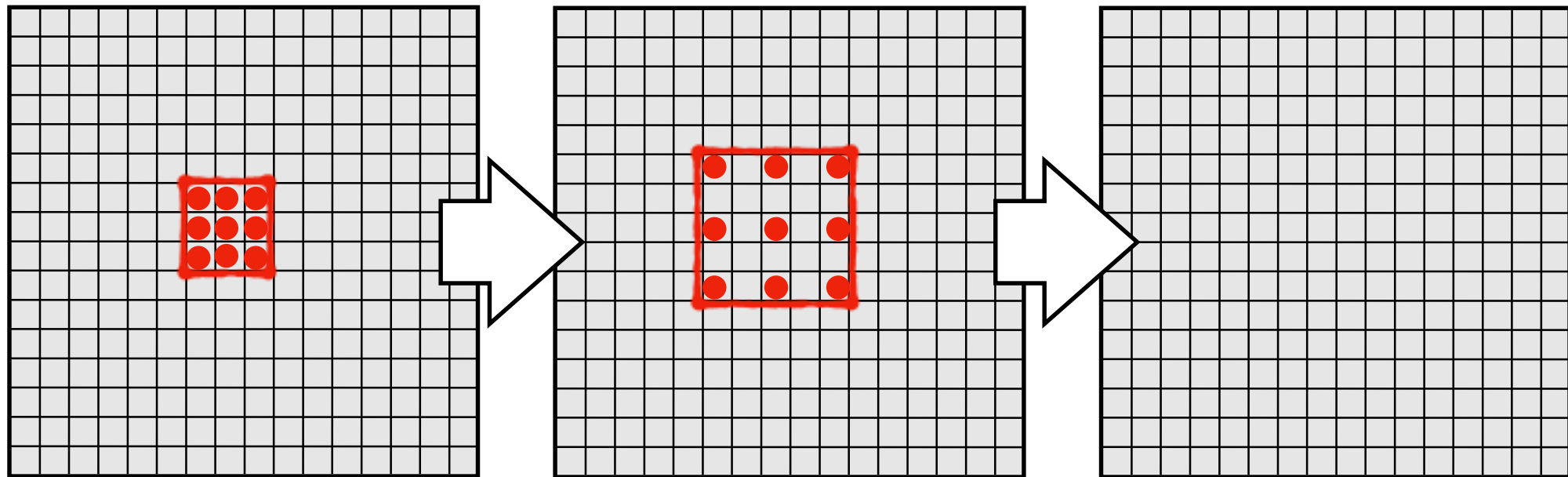
Receptive Field of Convolutions

Standard convolution:

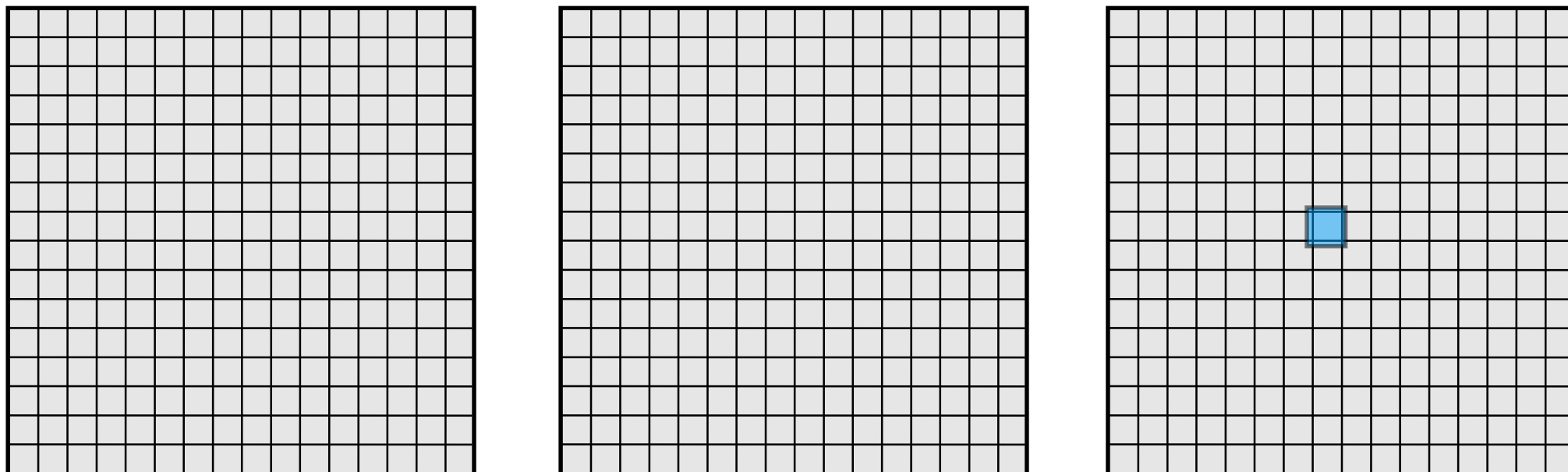


Receptive Field of Convolutions

Standard convolution:

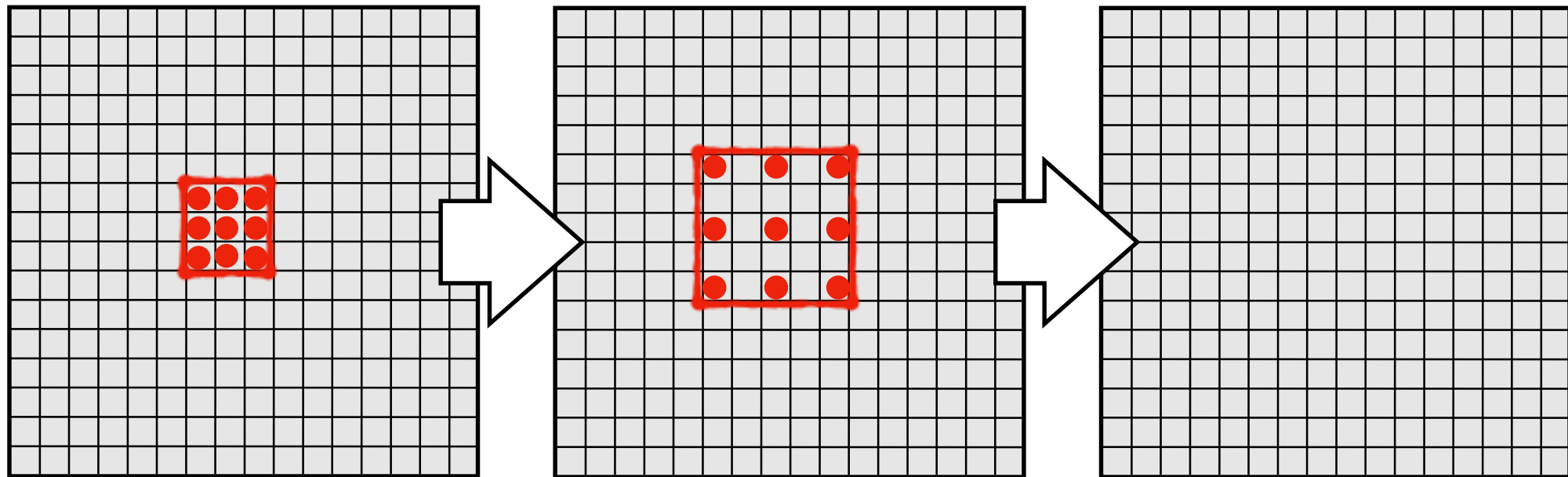


Receptive field:

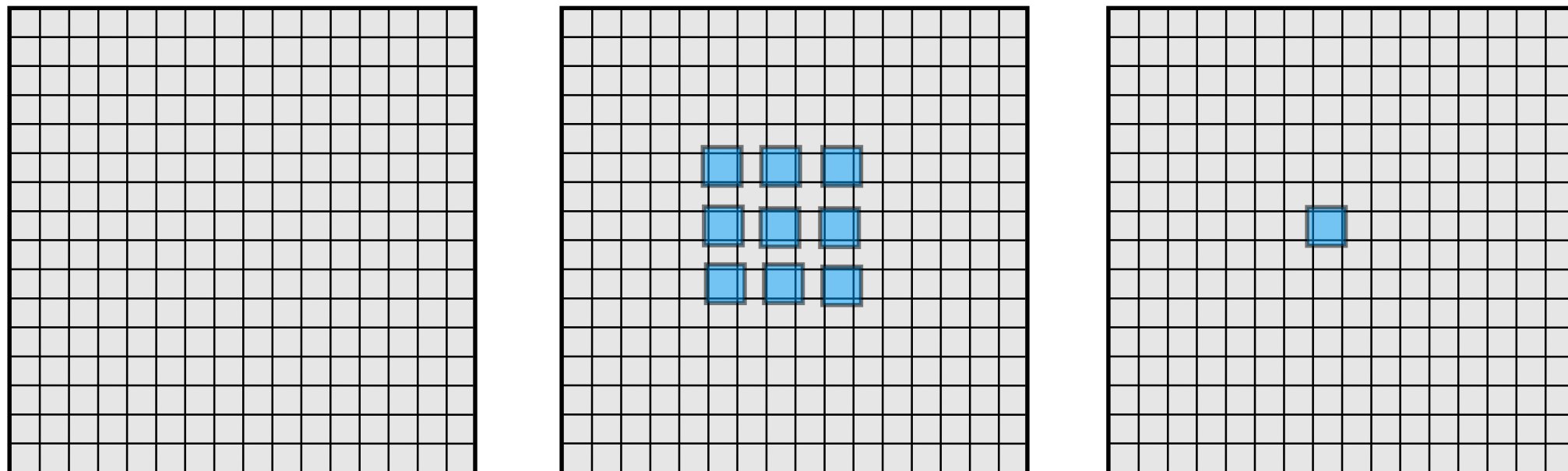


Receptive Field of Convolutions

Standard convolution:

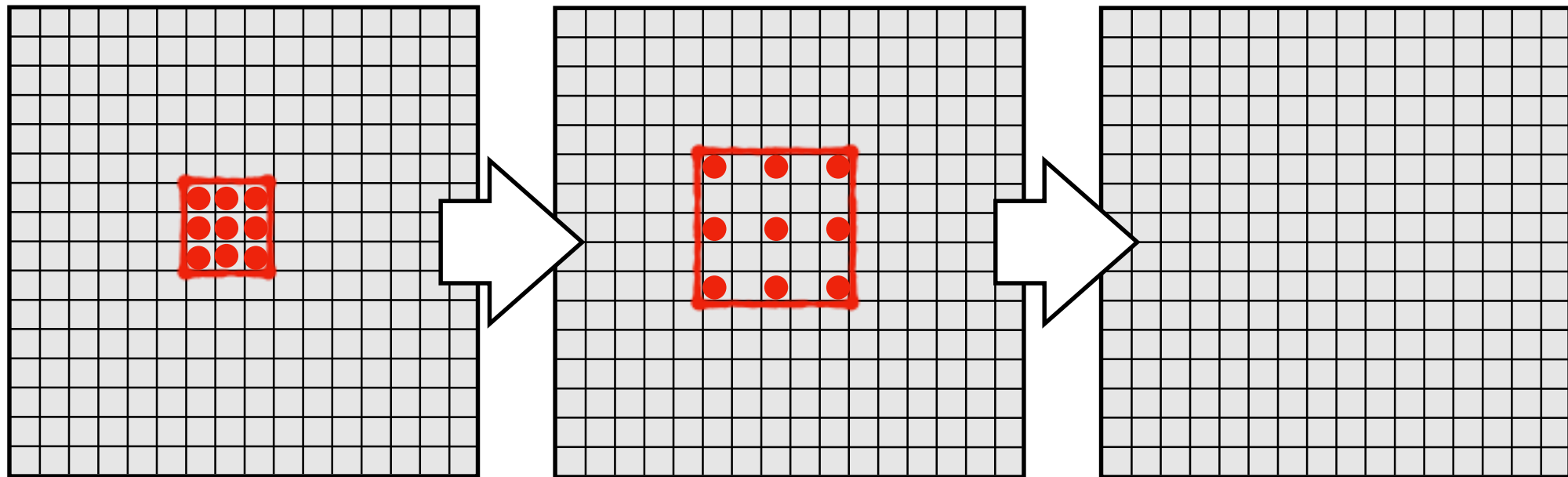


Receptive field:

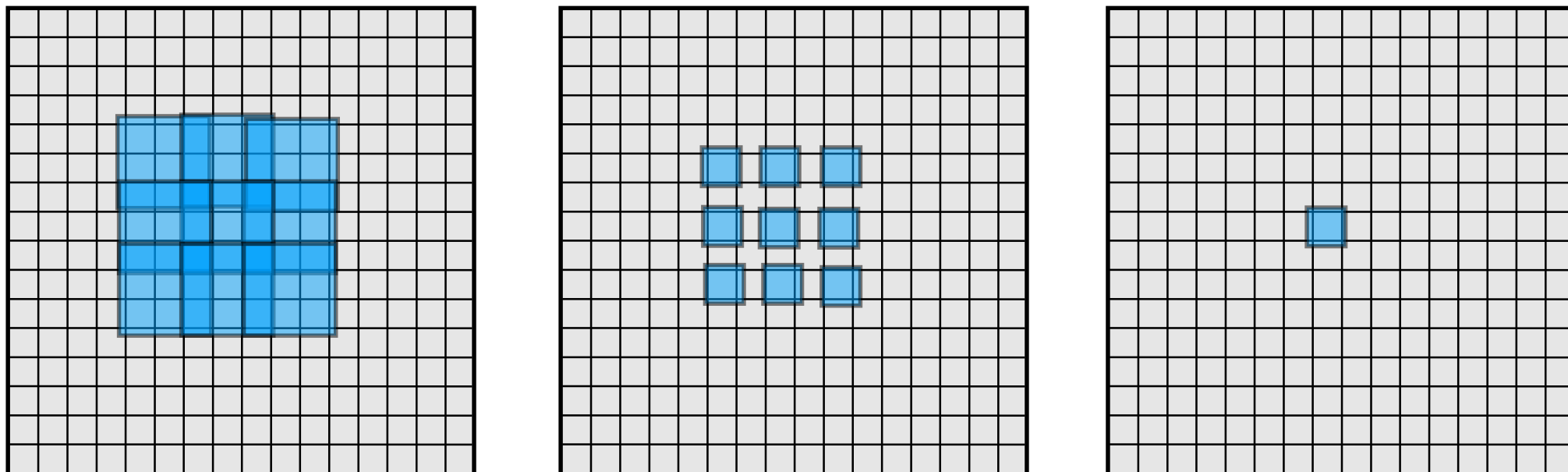


Receptive Field of Convolutions

Standard convolution:

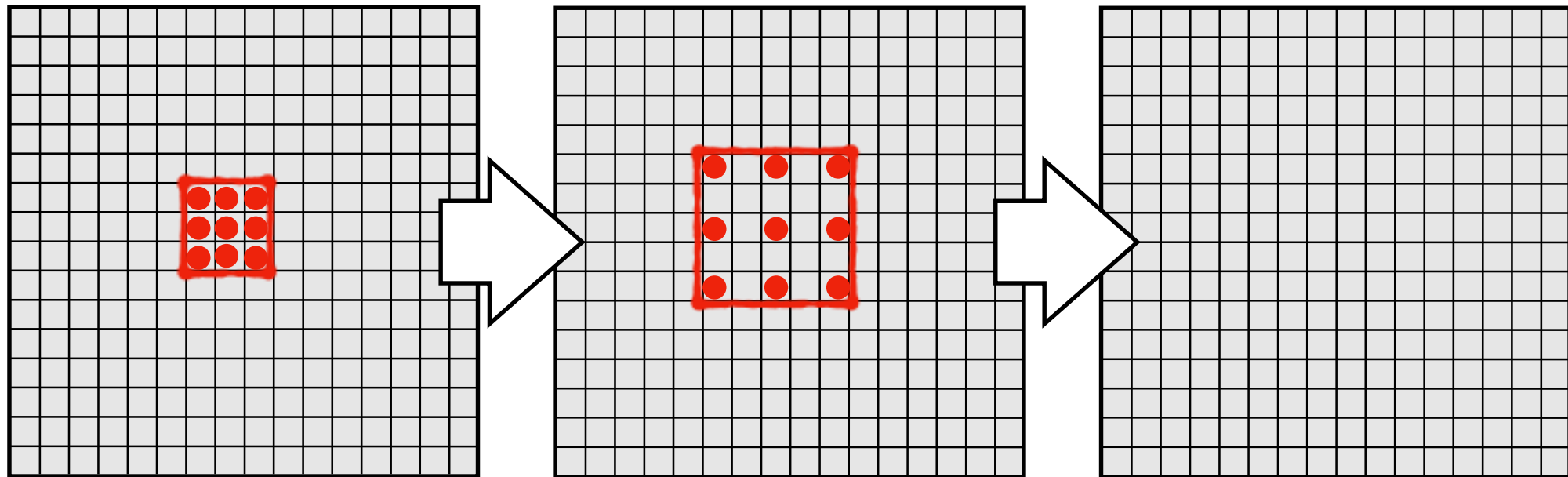


Receptive field:

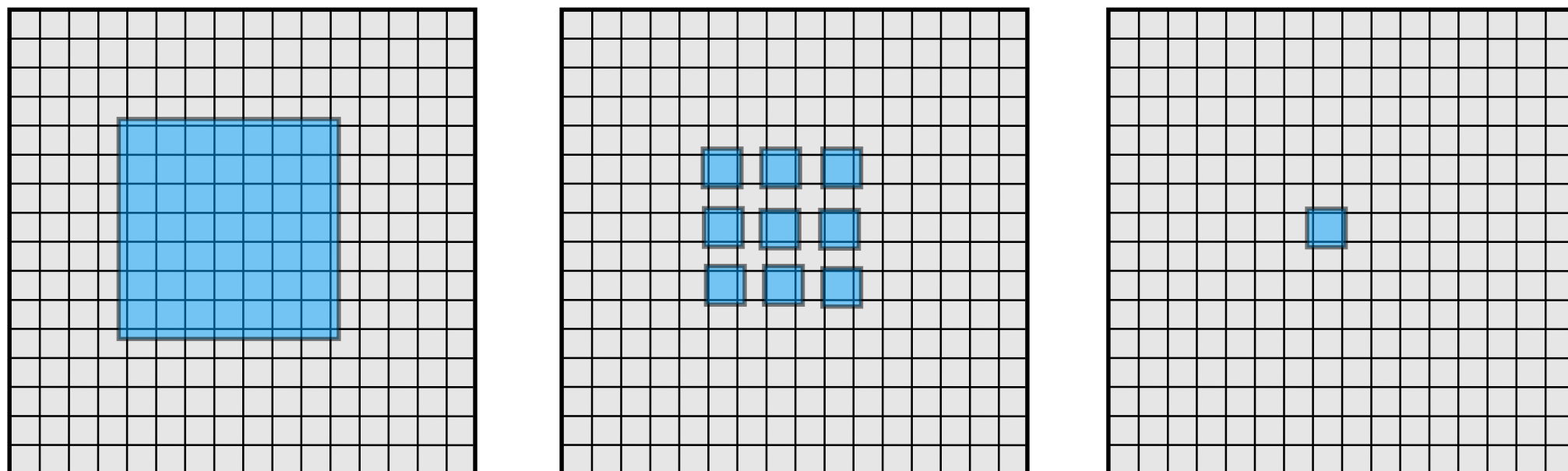


Receptive Field of Convolutions

Standard convolution:



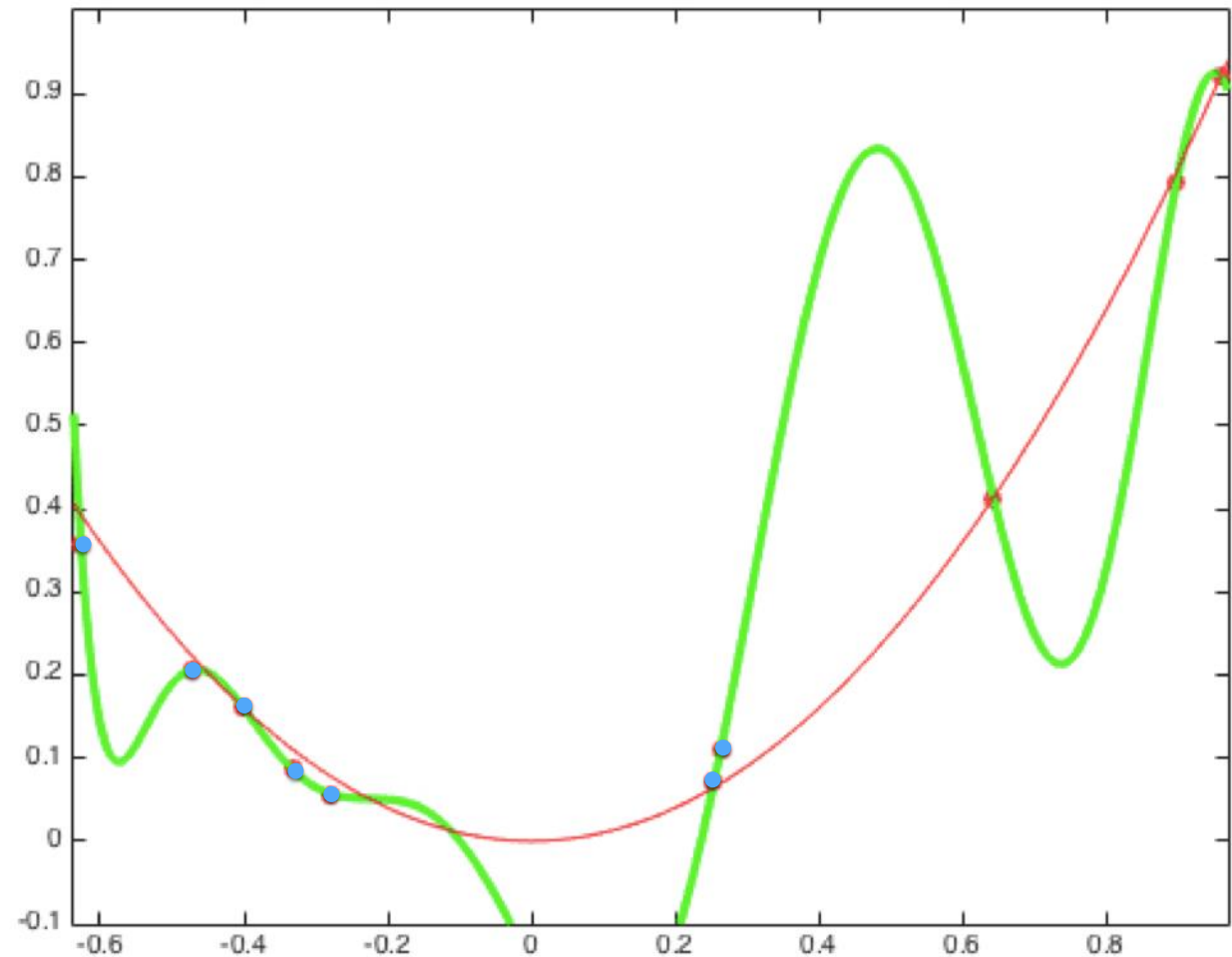
Receptive field:



Overfitting, Part II

Overfitting

10 data points - Fitting 10th degree polynomial



$$y = x^2$$

Overfitting

VGG-16 architecture

224



224

Overfitting

VGG-16 architecture



224

Conv
64 3x3

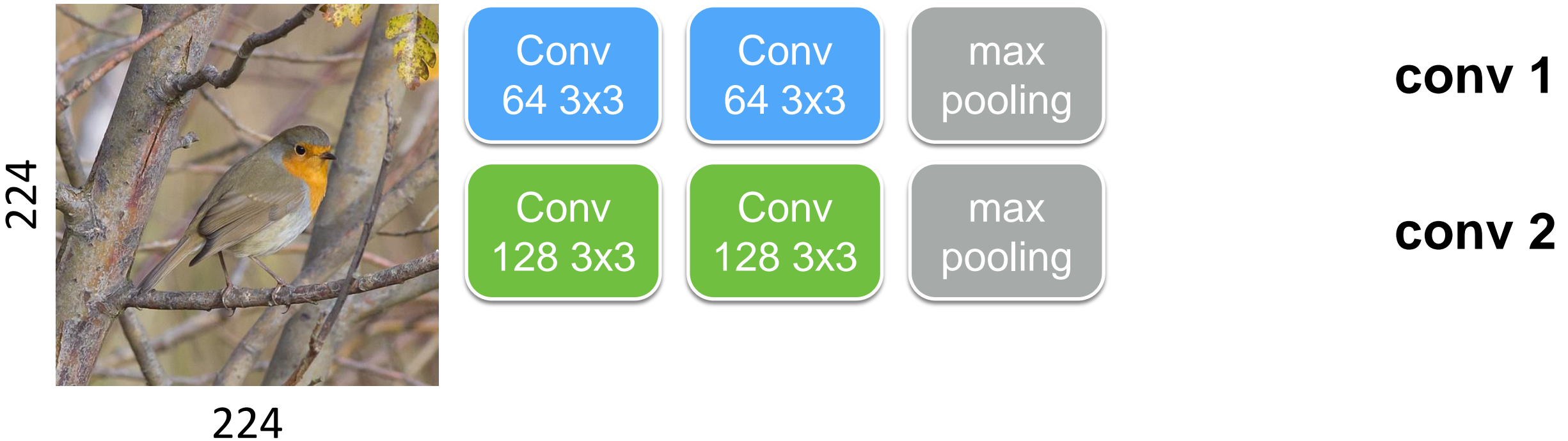
Conv
64 3x3

max
pooling

conv 1

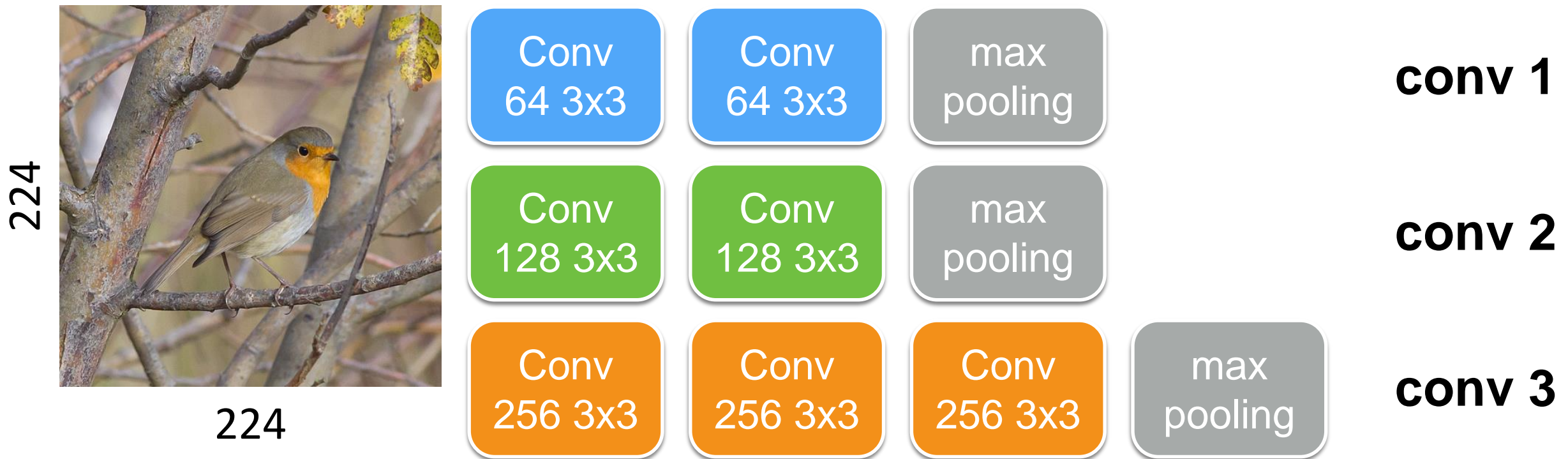
Overfitting

VGG-16 architecture



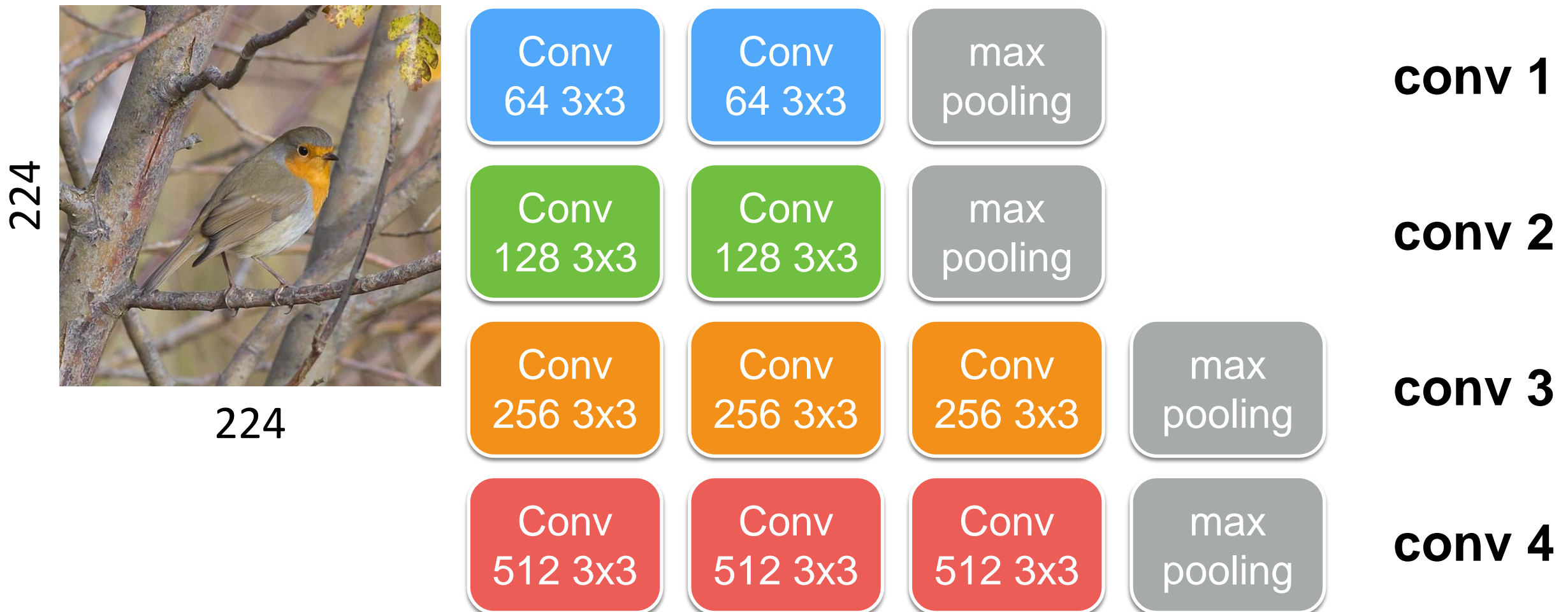
Overfitting

VGG-16 architecture



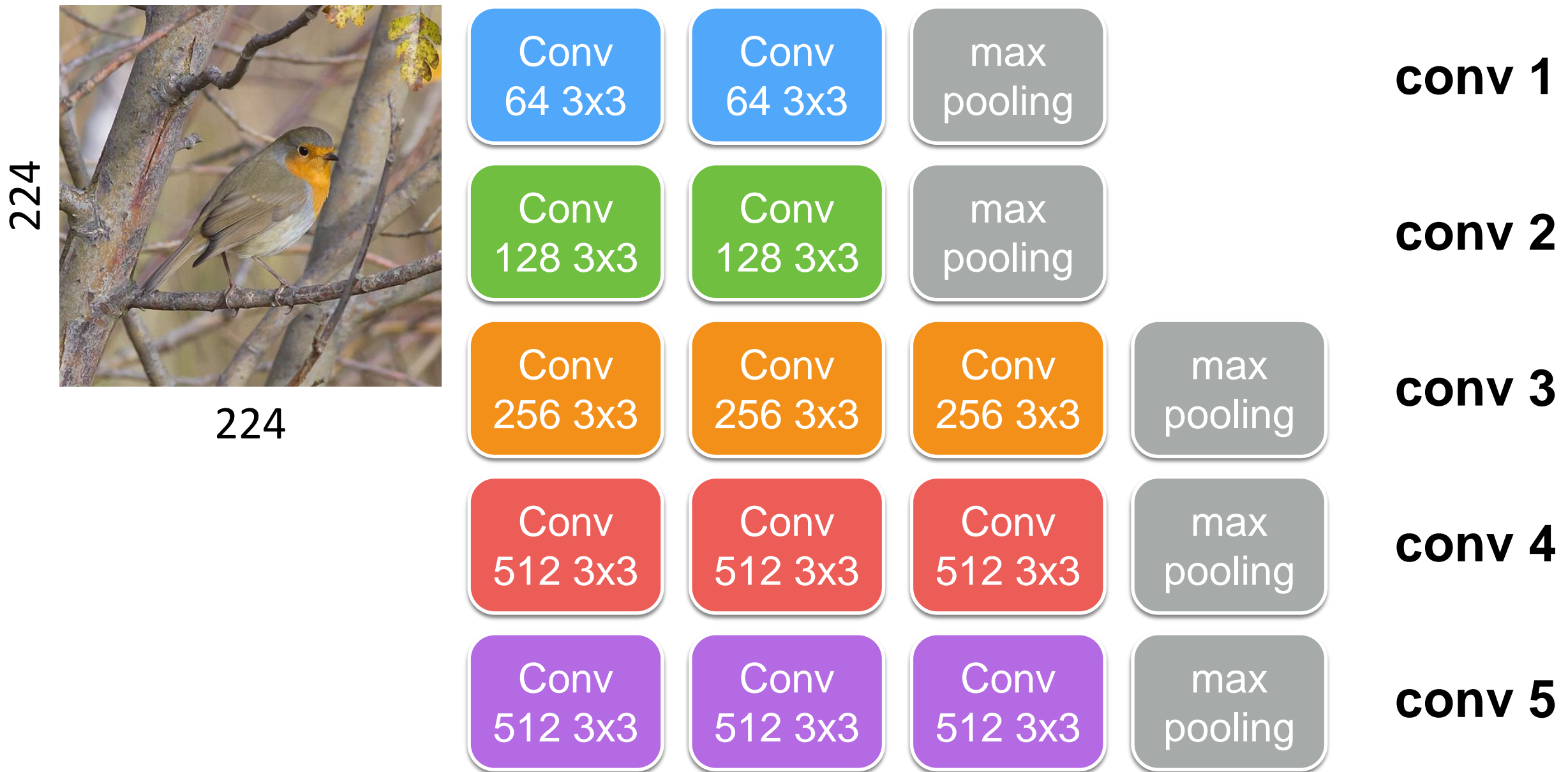
Overfitting

VGG-16 architecture



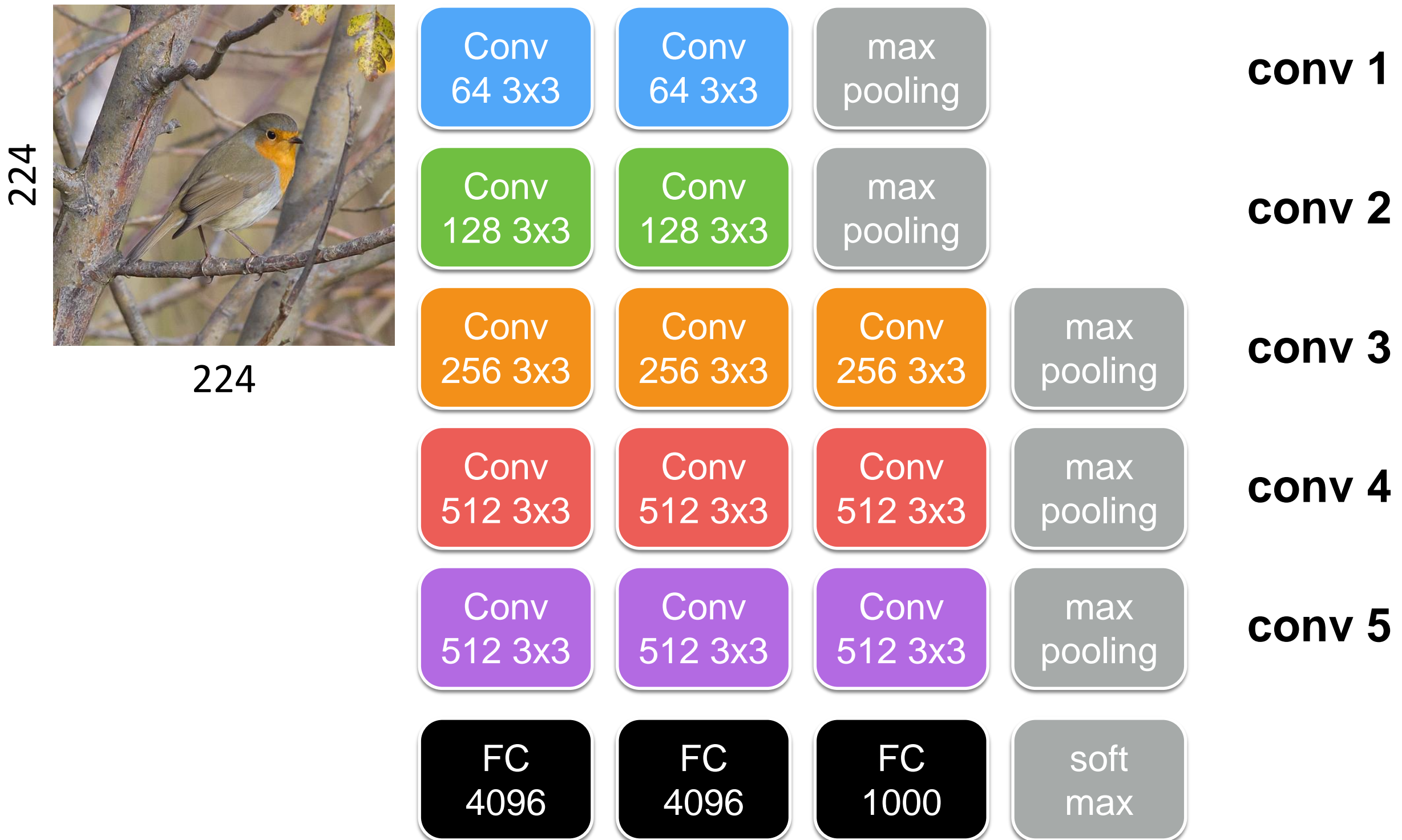
Overfitting

VGG-16 architecture



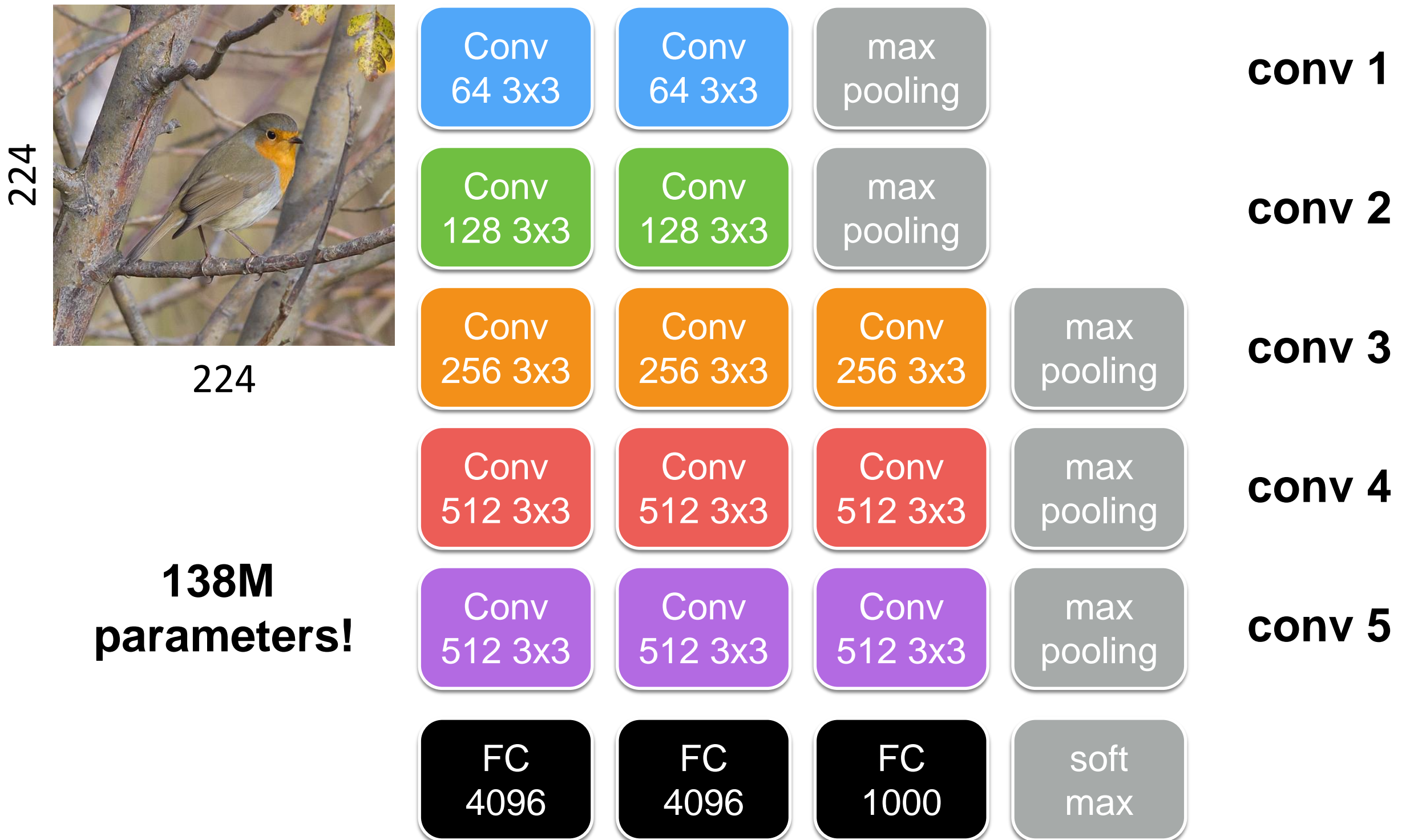
Overfitting

VGG-16 architecture

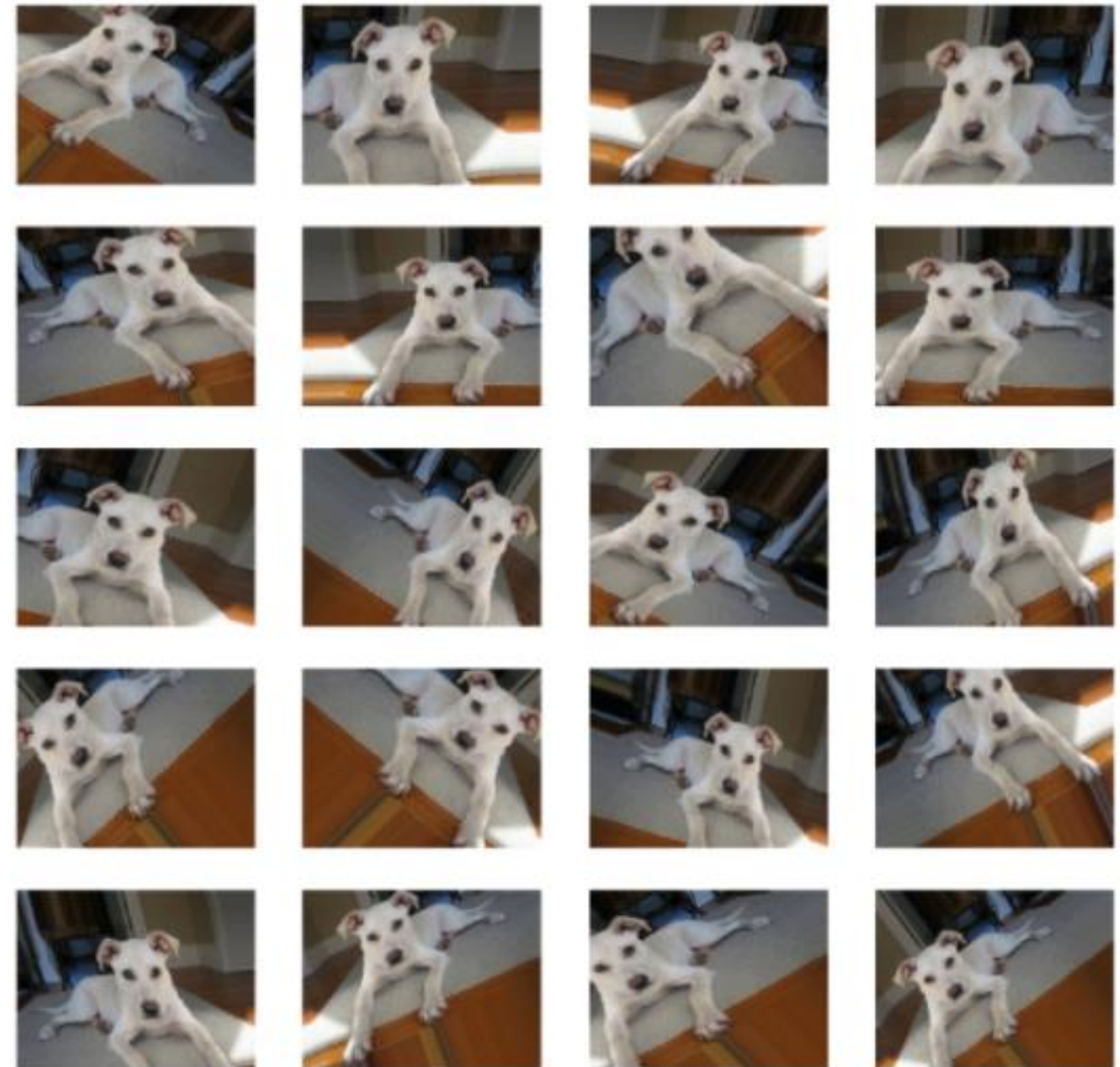


Overfitting

VGG-16 architecture



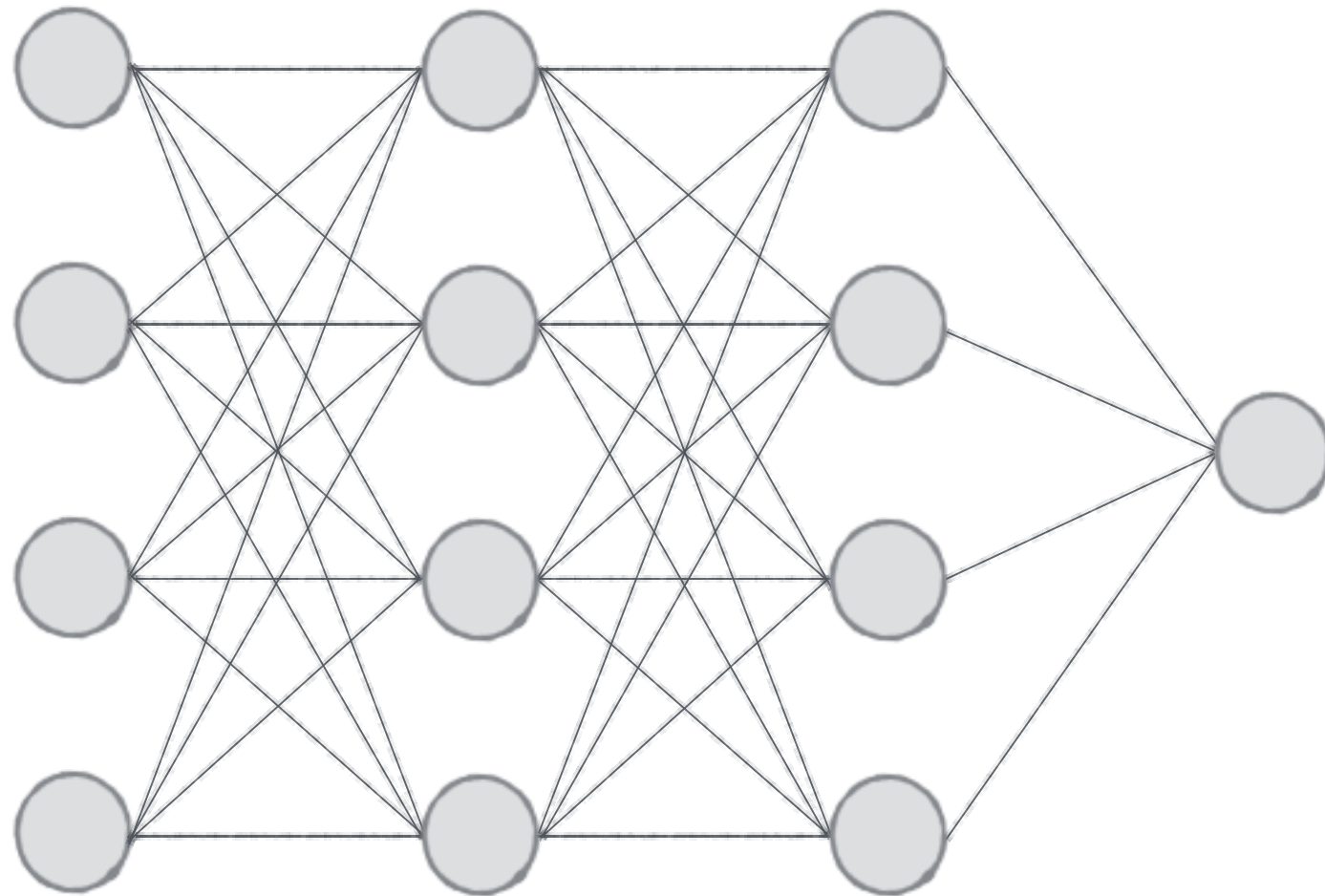
Data Augmentation



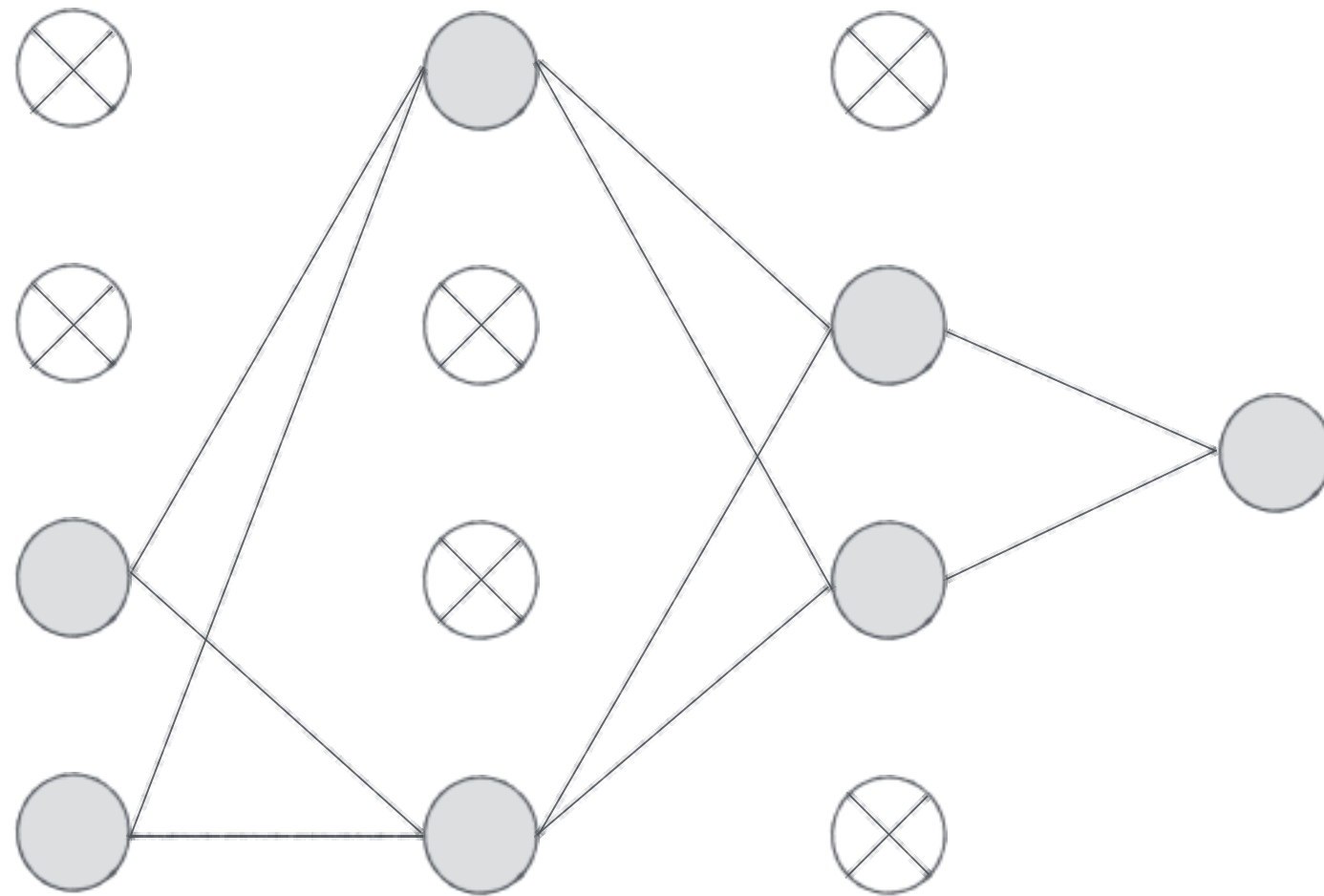
Regularization

$$L(w_1, \dots, w_n) + c \sum_i w_i^2$$

Dropout



Dropout



Randomly disable neurons for each minibatch

Dropout



Transfer Learning

IMAGENET



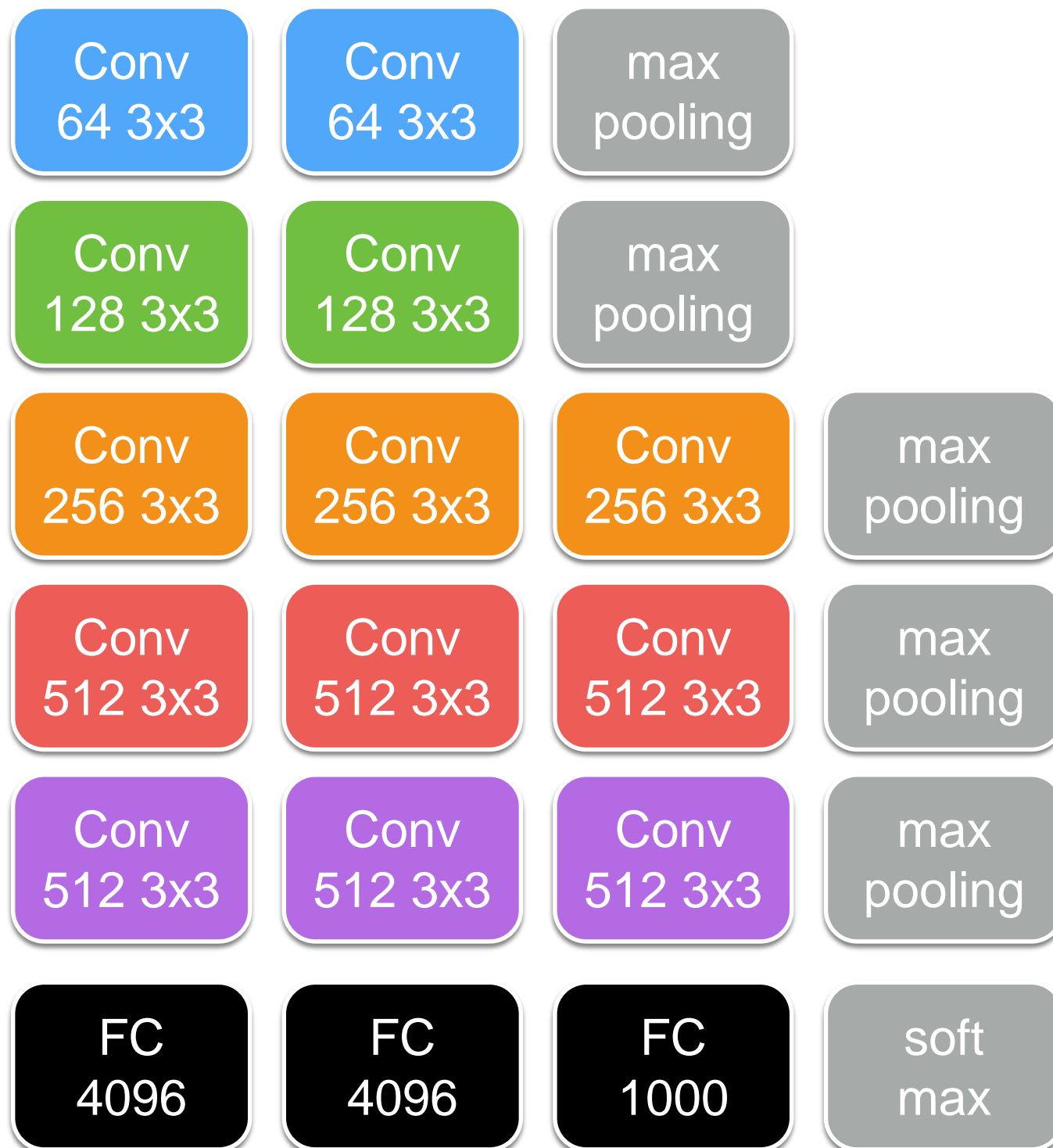
Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). [Imagenet large scale visual recognition challenge](#). *arXiv preprint arXiv:1409.0575*. [\[web\]](#)

3

14 million images of 1000 classes

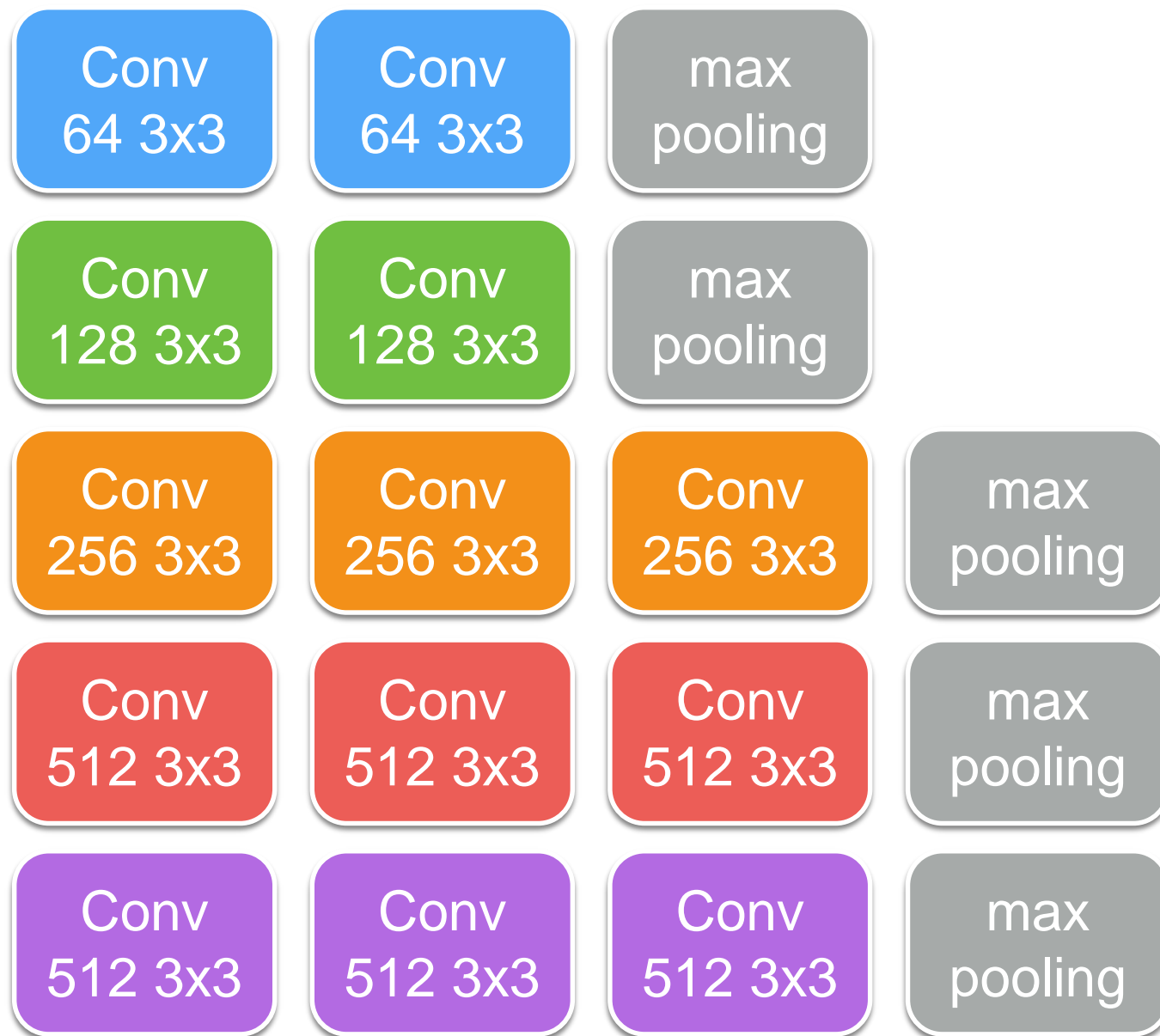
Transfer Learning

VGG-16 architecture



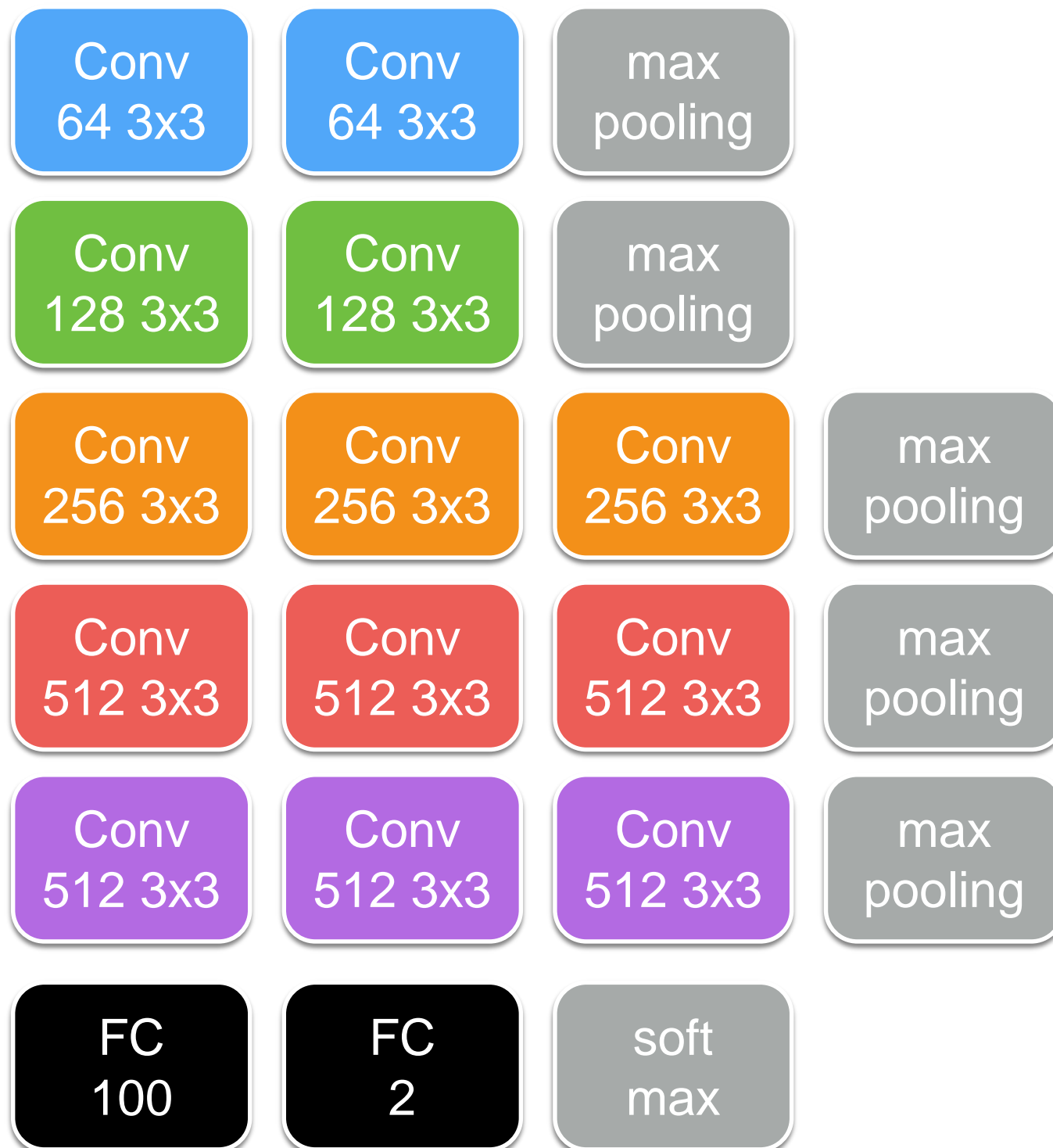
Transfer Learning

VGG-16 architecture



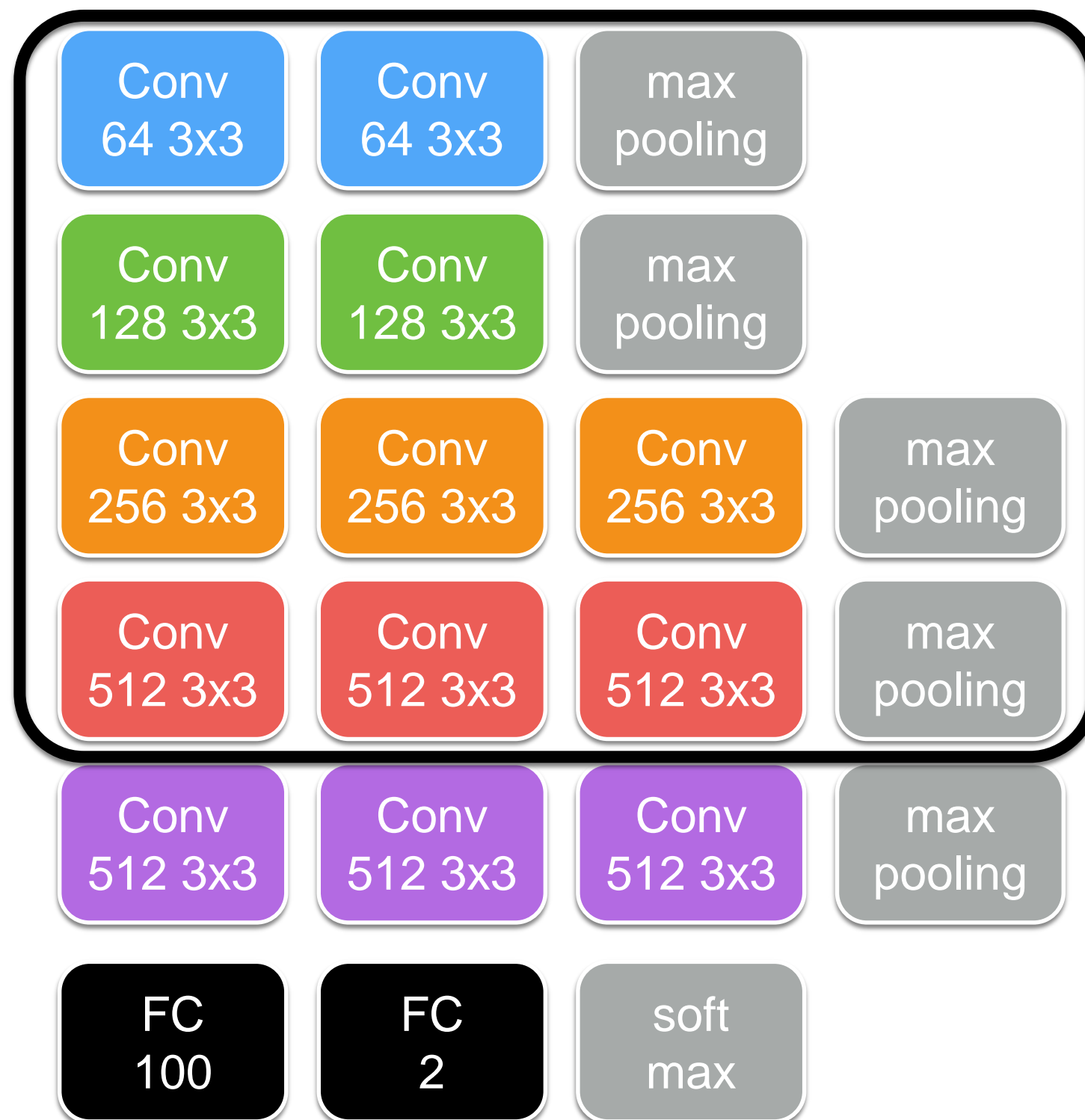
Transfer Learning

VGG-16 architecture



Transfer Learning

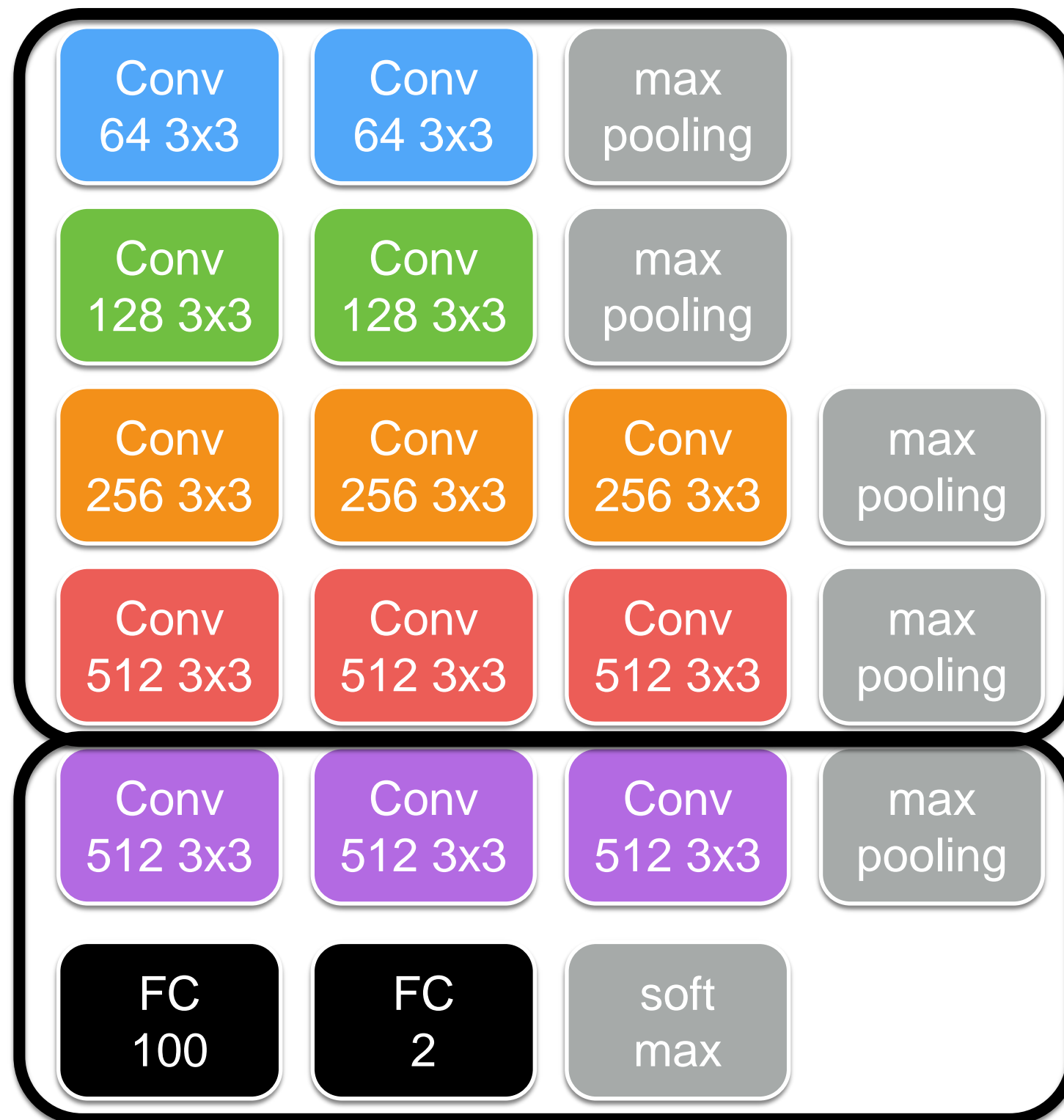
VGG-16 architecture



Fix weights

Transfer Learning

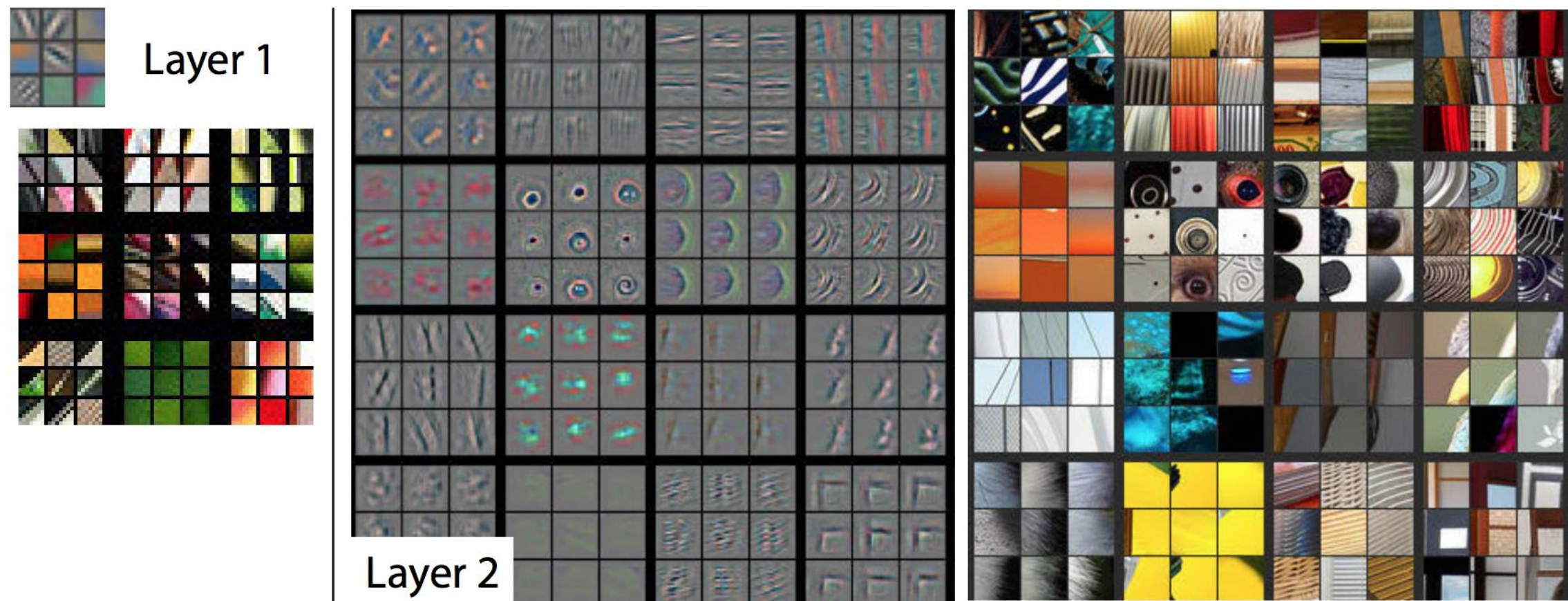
VGG-16 architecture



Fix weights

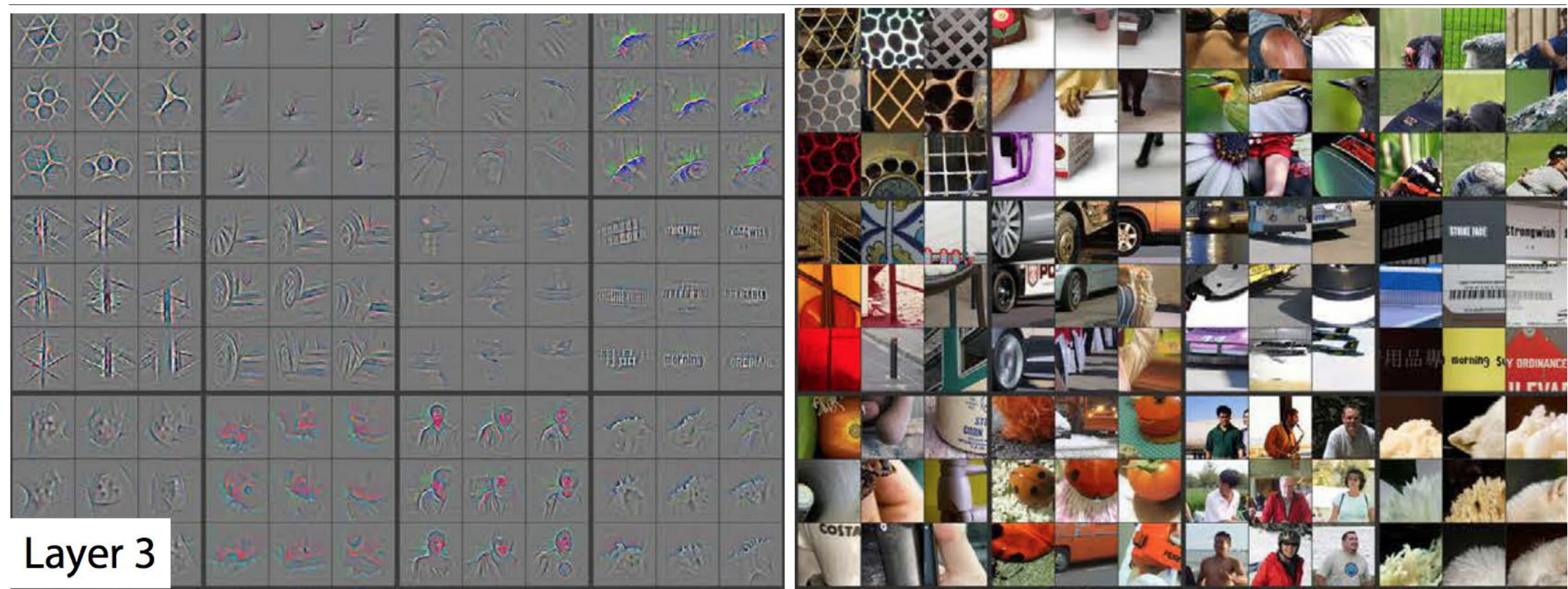
Train on limited data

Transfer Learning



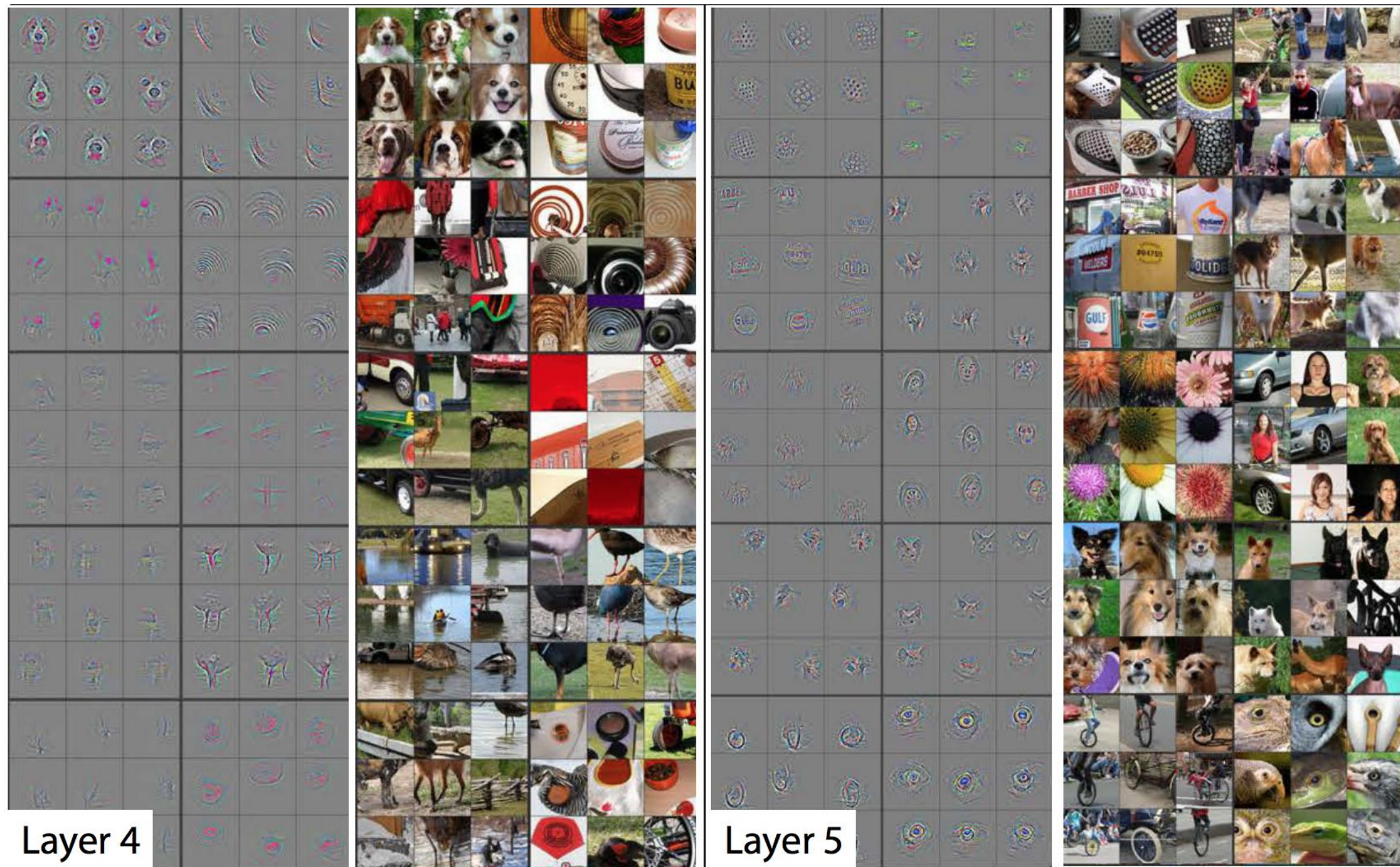
Low-level features: corners, edges, ...

Transfer Learning



Mid-level features

Transfer Learning



Higher-level features: object parts & objects

[Zeiler & Fergus, Visualizing and Understanding Convolutional Networks, ECCV 2014]

Lessons Learned

- Main lessons from this lecture
 - Convolutional neural networks
 - Layers: Fully connected, convolutional, activation functions, max pooling, strides & dilated filters
 - Fully convolutional networks
 - Handling overfitting
- Next lecture: Robust Model Fitting

Next Lecture

Jan. 20	Introduction, Linear classifiers and filtering	
Jan. 23	Filtering, gradients, scale	Lab 1
Jan. 27	Local features	
Jan. 30	Learning a classifier	Lab 2
Feb. 3	Convolutional neural networks	
Feb. 6	More convolutional neural networks	
Feb. 10	Robust model fitting and RANSAC	Lab 3
Feb. 13	Image registration	
Feb. 17	Camera Geometry	Lab 4
Feb. 20	More camera geometry	
Feb. 24	Generative neural networks	
Feb. 27	Generative neural networks	
Mar. 2	TBA	
Mar. 9	TBA	