



Week 4: Convolutional Networks

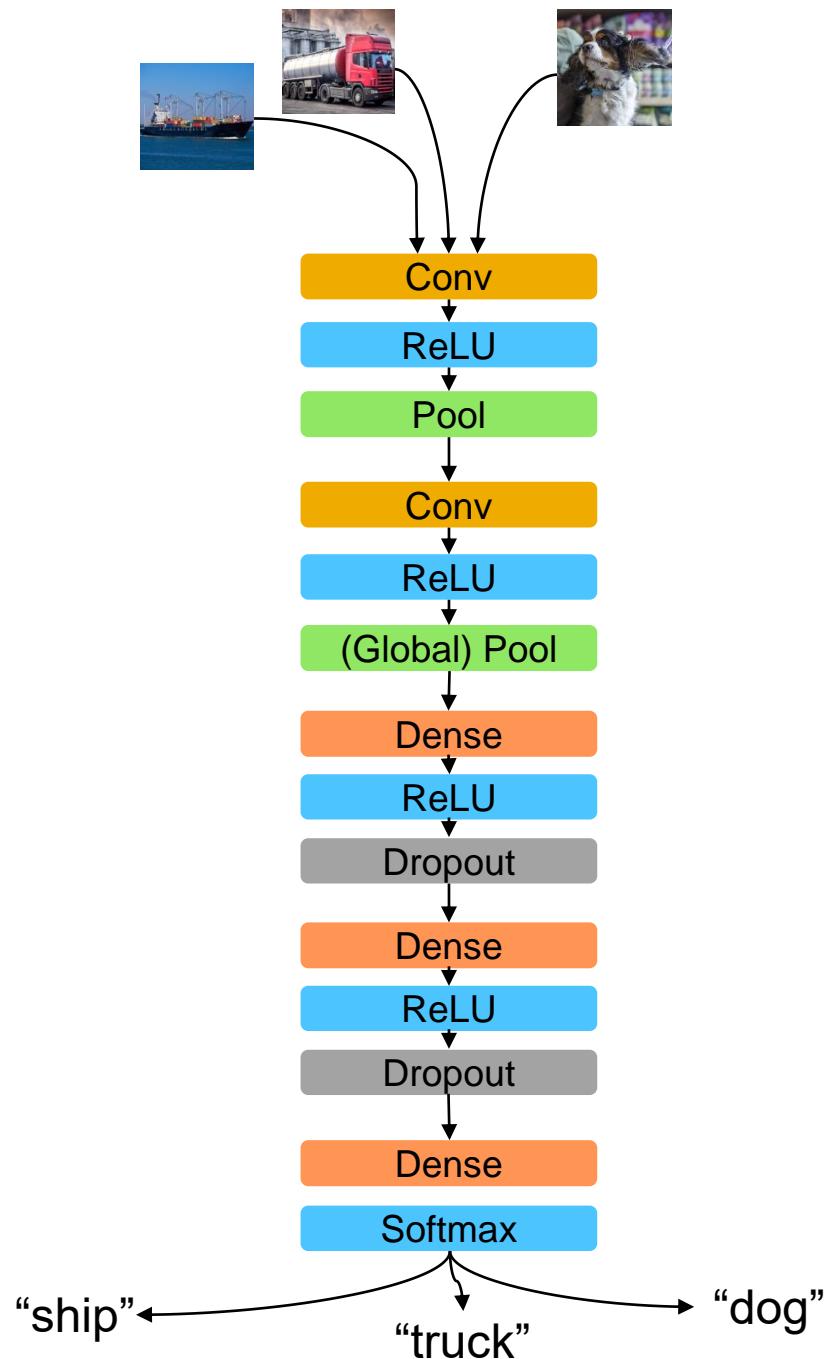
## **Unit 1: Introduction to CNNs**

# Introduction to CNNs

## Overview

### Content:

- Biological inspiration
- Challenges for computer vision
- The need for spatial invariance
- A naïve approach to neural networks for computer vision
- Convolutional neural networks

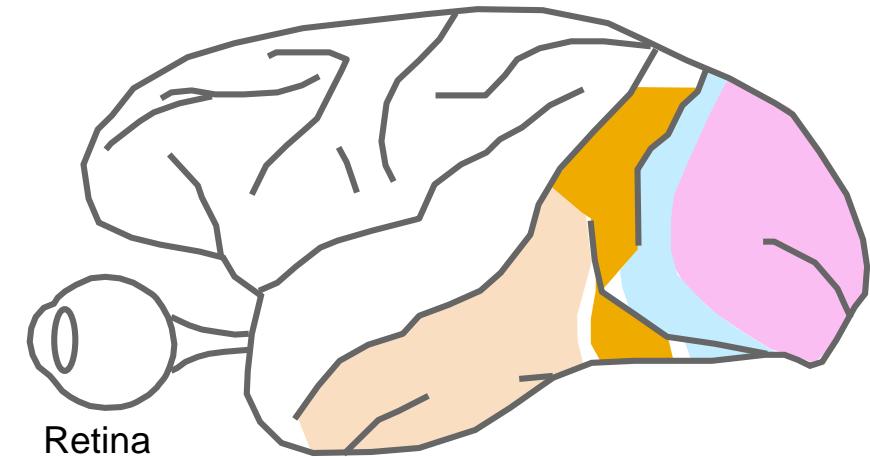


# Introduction to CNNs

Biological inspiration

## Example: The primate visual cortex

- Signals arriving from the retina are processed hierarchically by subsequent brain areas
- This is reminiscent of processing by subsequent layers of a deep artificial neural network



David Daniel Cox, Thomas Dean, *Neural Networks and Neuroscience-Inspired Computer Vision*, Current Biology, Volume 24, Issue 18, 2014, Pages R921-R929

# Introduction to CNNs

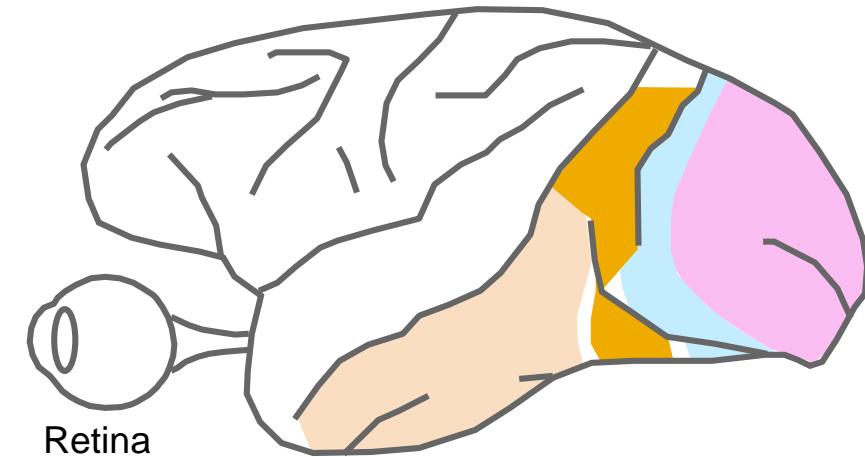
Biological inspiration

## Example: The primate visual cortex

- Signals arriving from the retina are processed hierarchically by subsequent brain areas
- This is reminiscent of processing by subsequent layers of a deep artificial neural network

**Of course, the correspondence is not perfect:**

- In the primate visual cortex, there are many forward and backward connections



David Daniel Cox, Thomas Dean, *Neural Networks and Neuroscience-Inspired Computer Vision*, Current Biology, Volume 24, Issue 18, 2014, Pages R921-R929

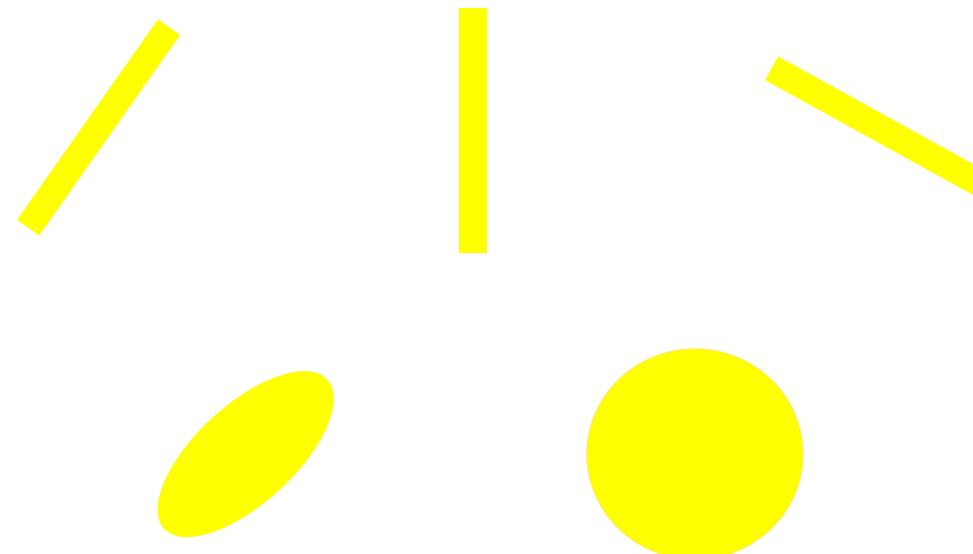
# Introduction to CNNs

Biological inspiration

Hubel and Wiesel (1950s and '60s)  
studied feline visual cortex

**Two types of cells identified:**

- *Simple* cells fire in response to stimuli of a particular shape and orientation



Hubel, David H., and Torsten N. Wiesel. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex." *The Journal of physiology* 160.1 (1962): 106

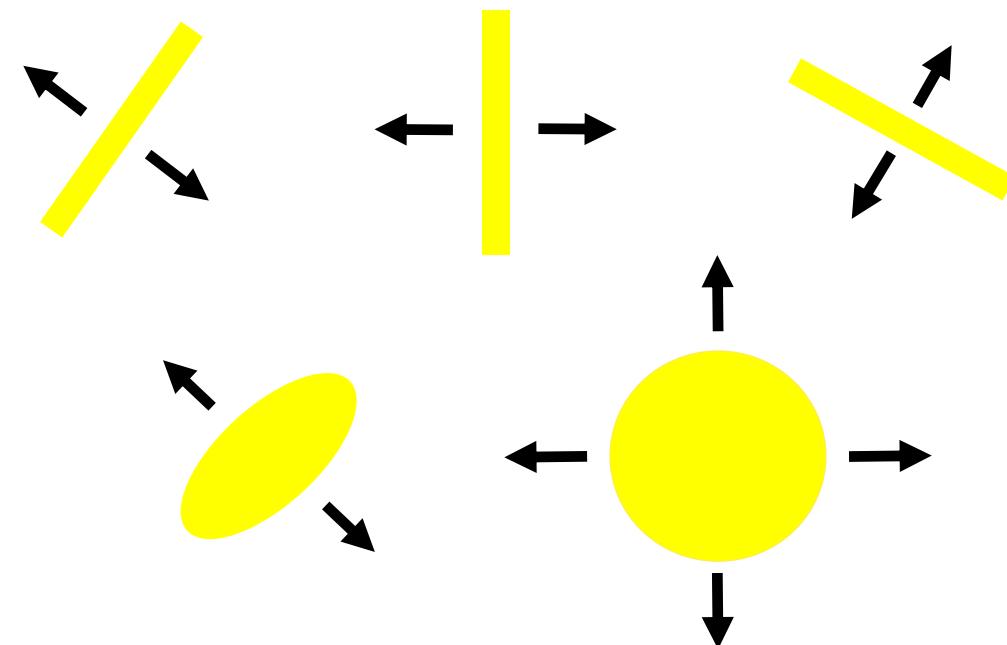
# Introduction to CNNs

Biological inspiration

Hubel and Wiesel (1950s and '60s)  
studied feline visual cortex

**Two types of cells identified:**

- *Simple* cells fire in response to stimuli of a particular shape and orientation
- *Complex* cells additionally fire only when the stimulus moves in a particular direction



Hubel, David H., and Torsten N. Wiesel. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex." *The Journal of physiology* 160.1 (1962): 106

# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye
- Identify other body parts of the dog, like a second eye



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye
- Identify other body parts of the dog, like a second eye, ears



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye
- Identify other body parts of the dog, like a second eye, ears, snout



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye
- Identify other body parts of the dog, like a second eye, ears, snout, legs, etc.



# Introduction to CNNs

Challenges for computer vision – From images to abstract representations

Object recognition involves combining many irregular features into a whole

**For example, to recognize a dog, an animal's brain must:**

- Identify edges of shapes, like the arcs forming the outline of an eye
- Combine those component parts into an abstract representation of an eye
- Identify other body parts of the dog, like a second eye, ears, snout, legs, etc.
- Combine these parts into an abstract representation of a dog



# Introduction to CNNs

Challenges for computer vision – Data dimensionality

Image data is very high-dimensional

- A  $512 \times 512$  RGB image has  $512 \cdot 512 \cdot 3 = 786,432$  features

# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*



# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*

Many invariances must be learned by a “naïve” machine learning model, e.g.:

- the position of the dog in the image (spatial invariance)



# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*

Many invariances must be learned by a “naïve” machine learning model, e.g.:

- the position of the dog in the image (spatial invariance)
- the size of the dog in the image (scale invariance)



# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*

Many invariances must be learned by a “naïve” machine learning model, e.g.:

- the position of the dog in the image (spatial invariance)
- the size of the dog in the image (scale invariance)
- the angle at which the image was taken (2D rotational invariance)



# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*

Many invariances must be learned by a “naïve” machine learning model, e.g.:

- the position of the dog in the image (spatial invariance)
- the size of the dog in the image (scale invariance)
- the angle at which the image was taken (2D rotational invariance)
- the orientation of the dog and the parts of its body in 3D space



# Introduction to CNNs

Challenges for computer vision – Invariance

*Does this image contain a dog?*

Many invariances must be learned by a “naïve” machine learning model, e.g.:

- the position of the dog in the image (spatial invariance)
- the size of the dog in the image (scale invariance)
- the angle at which the image was taken (2D rotational invariance)
- the orientation of the dog and the parts of its body in 3D space
- and more...



# Introduction to CNNs

The need for spatial invariance

*Does this image contain a dog?*

Let's focus on spatial invariance.

- Where the dog appears is irrelevant
- The same is true for the component questions:
  - *Is there an edge in some orientation?*
  - *Is this an eye?*
  - *Does it have legs?*



# Introduction to CNNs

The need for spatial invariance

*Does this image contain a dog?*

Let's focus on spatial invariance.

- Where the dog appears is irrelevant
- The same is true for the component questions:
  - *Is there an edge in some orientation?*
  - *Is this an eye?*
  - *Does it have legs?*
- *However, the spatial relationships between the components of the image remain important*



dog

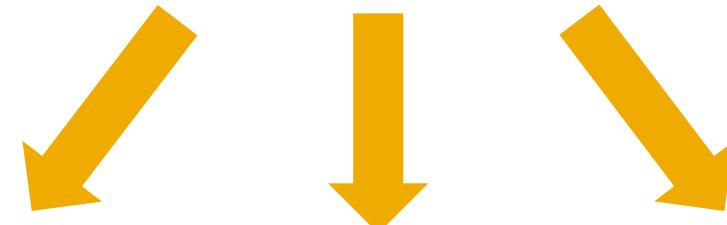


not a dog

# Introduction to CNNs

A naïve approach to neural networks for computer vision

A first idea for dealing with spatial invariance:  
Augment the data by generating many  
translations of the images in the training set,  
then apply a feed-forward neural network.



# Introduction to CNNs

A naïve approach to neural networks for computer vision

A first idea for dealing with spatial invariance:

Augment the data by generating many translations of the images in the training set, then apply a feed-forward neural network.

## Problems:

- Lengthens training time by massively increasing size of data set
- The feed-forward network requires many parameters (for a  $512 \times 512$  RGB image, 786,432 times the size of the first hidden layer in the input layer alone)



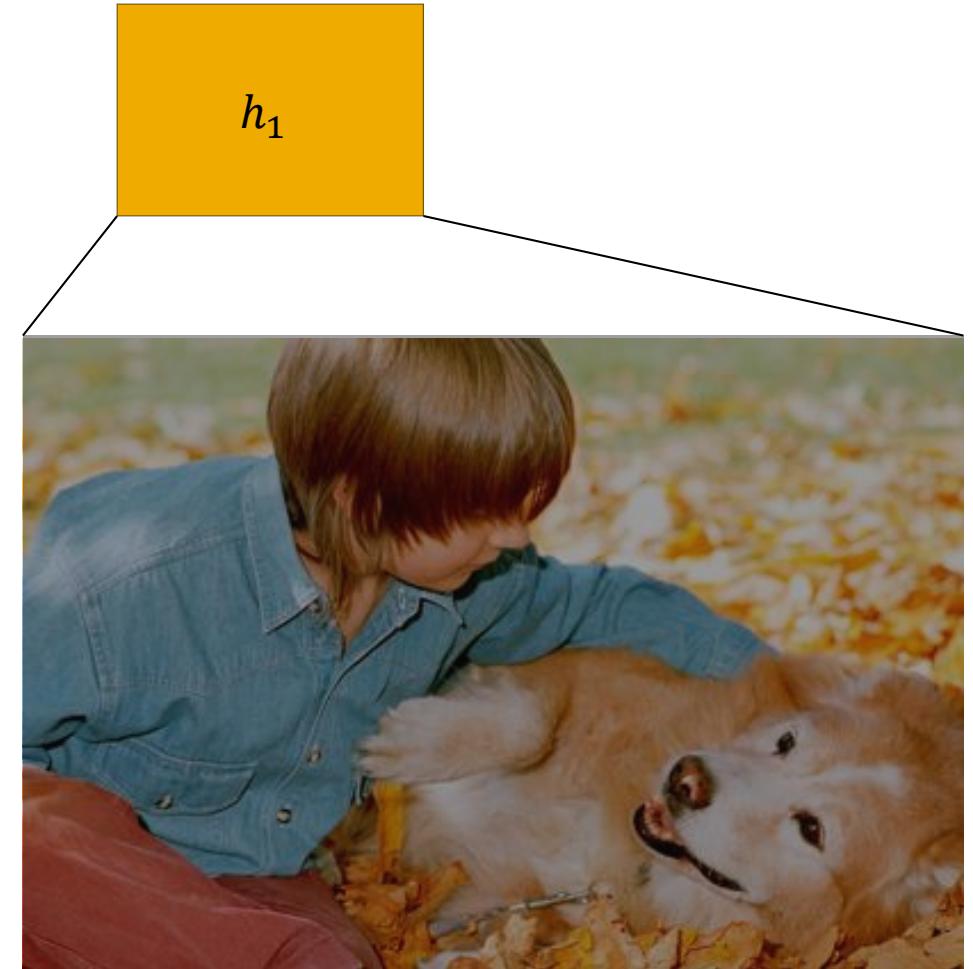
# Introduction to CNNs

A naïve approach to neural networks for computer vision

A feed-forward neural network connects each pixel to each hidden node

**A lot of parameters:**

$image\ width \times image\ height \times size\ of\ hidden\ layer$



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

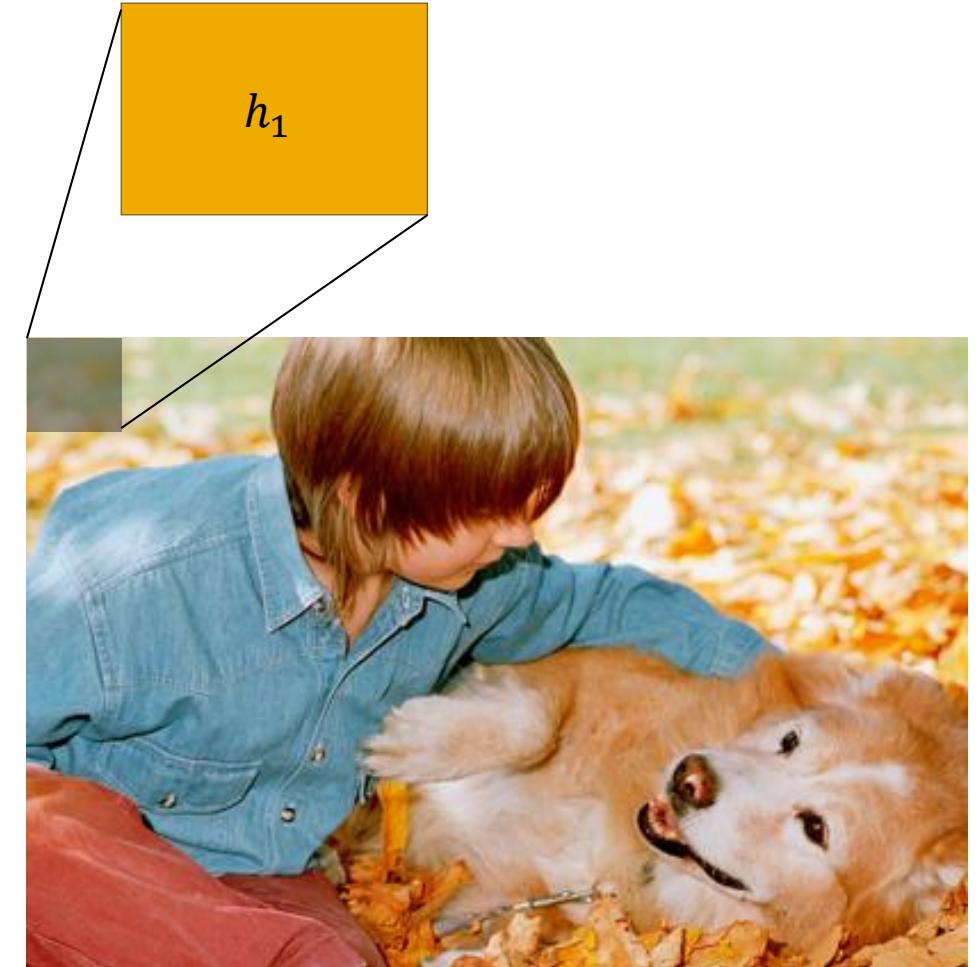


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

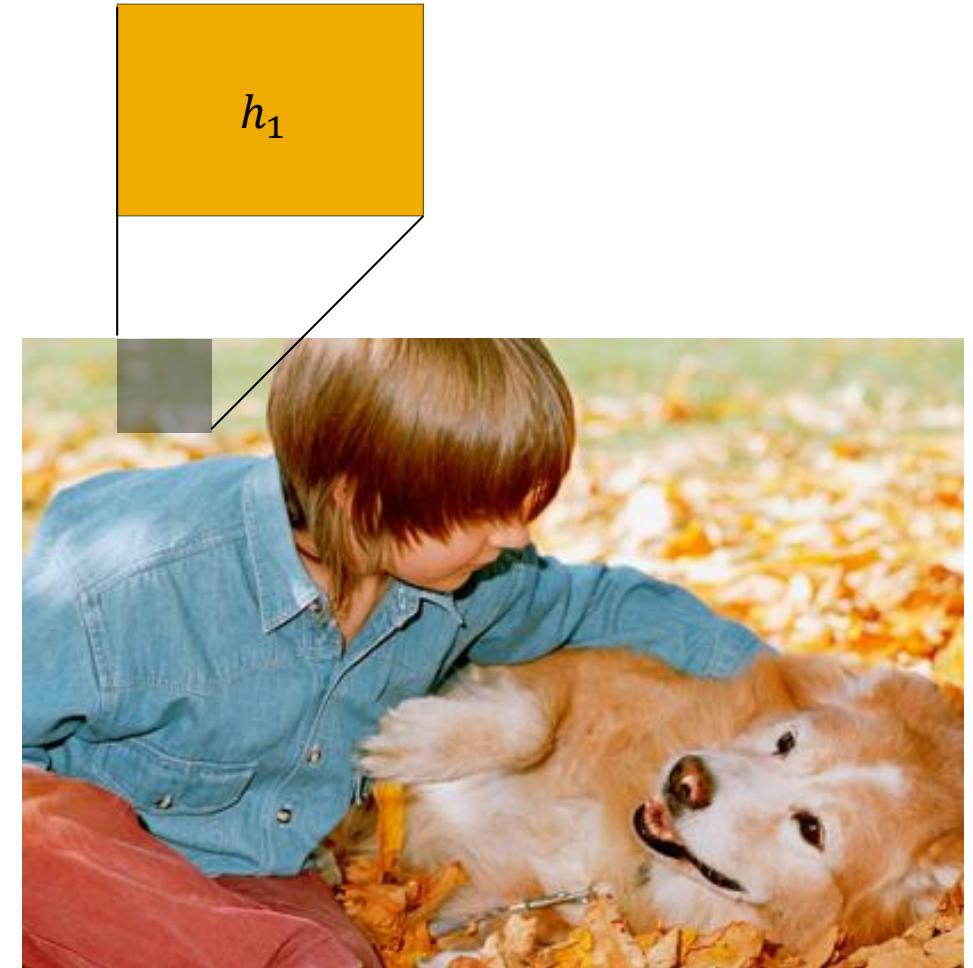


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

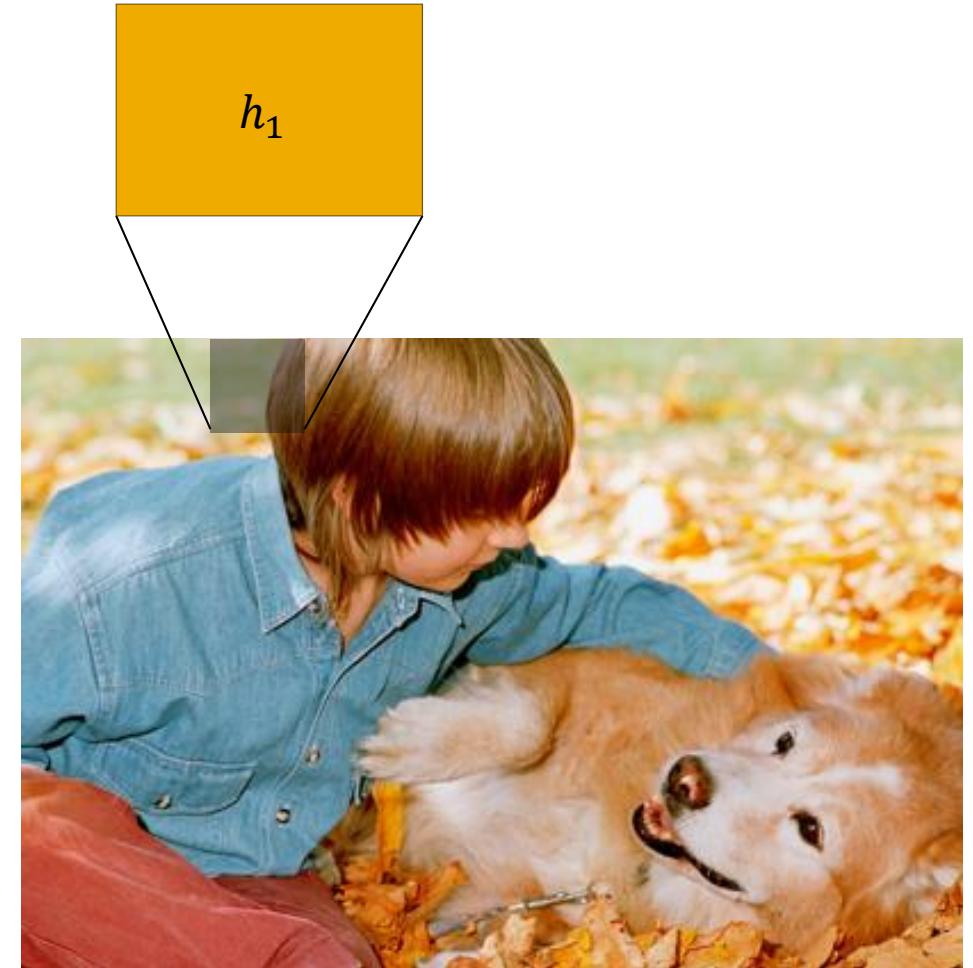


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

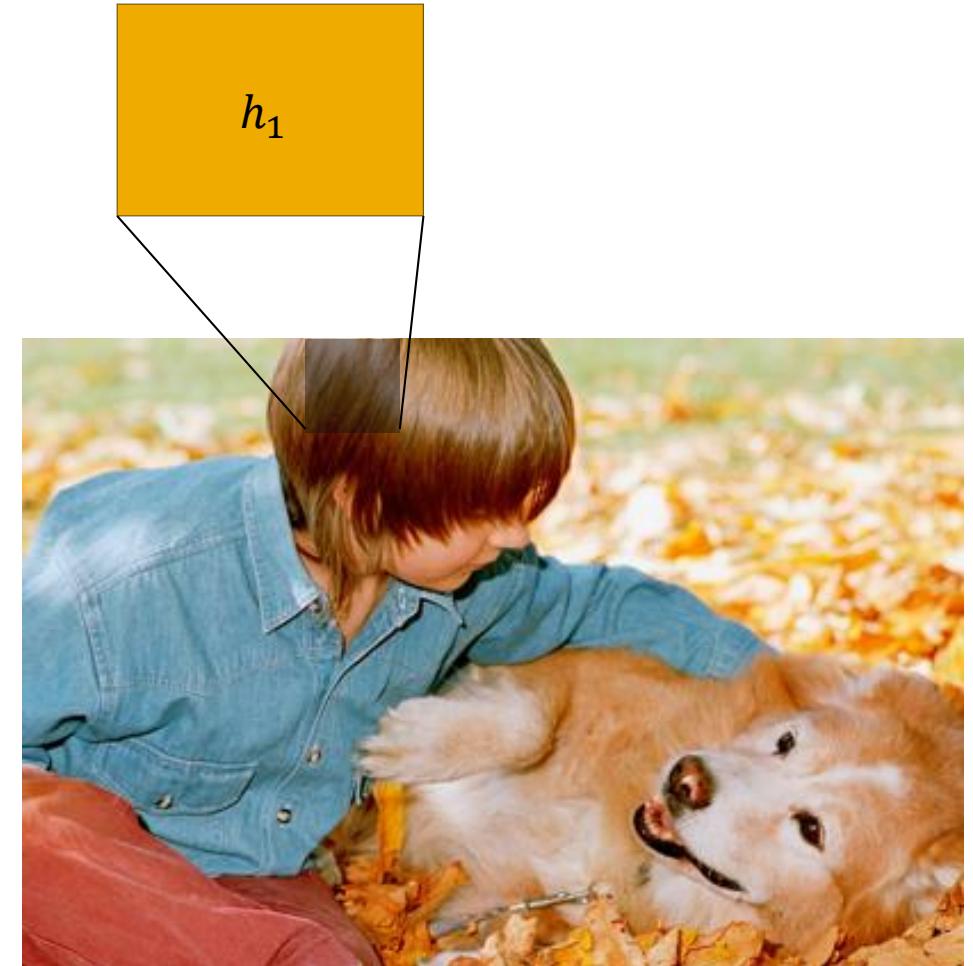


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

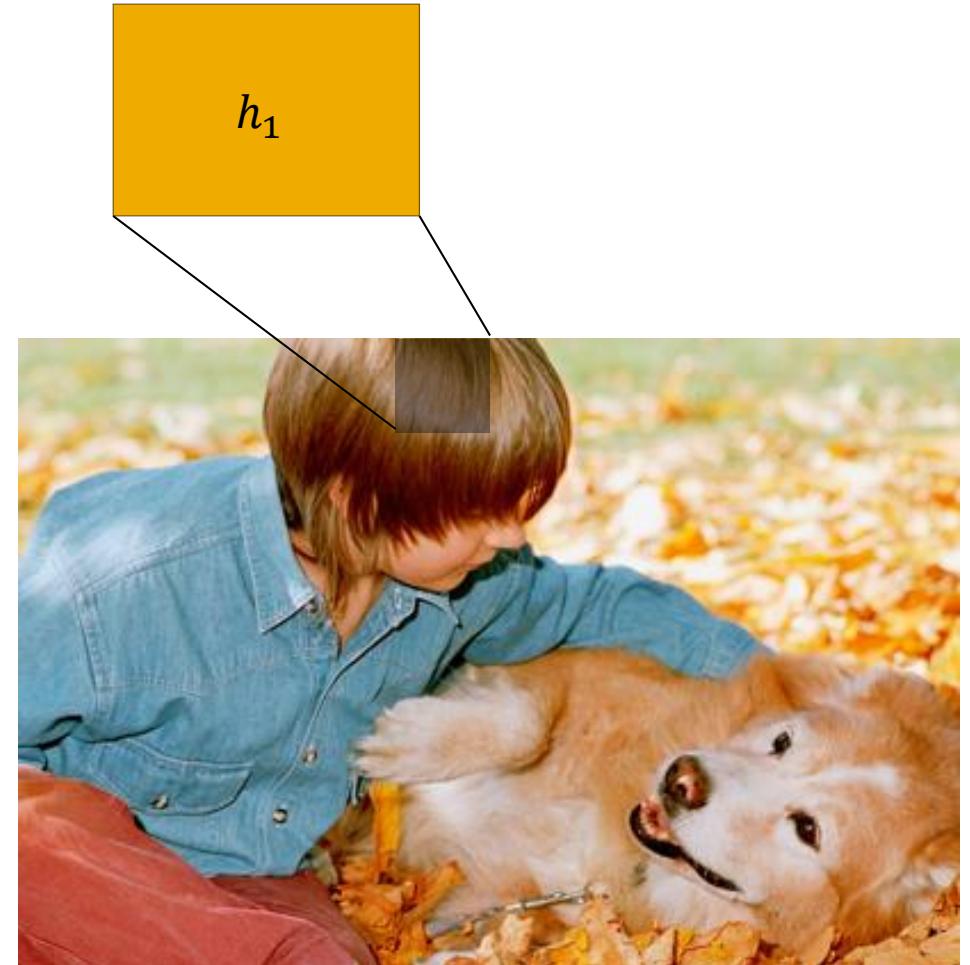


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

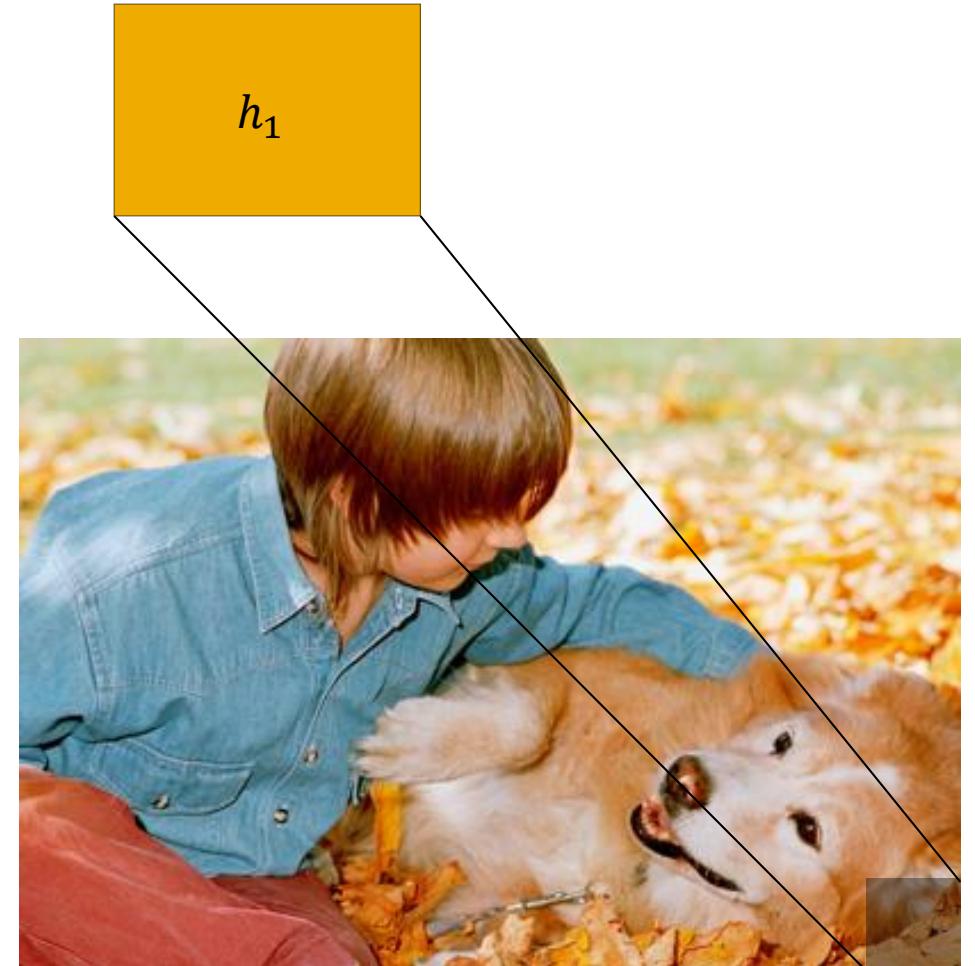


# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network



# Introduction to CNNs

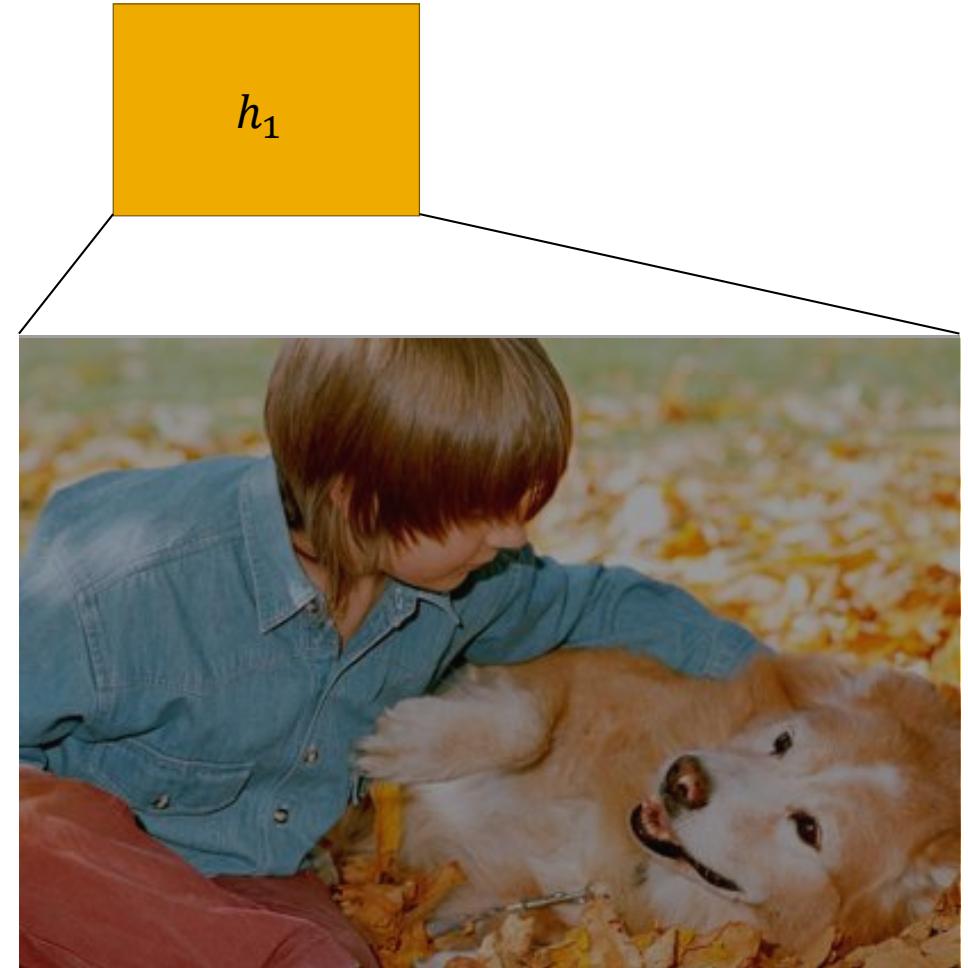
Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

Each small patch is connected to the first hidden layer of the network

Contrast with a feed-forward neural network, which connects each pixel to each hidden node

Significantly more parameters:  
*image width × image height × size of hidden layer*



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

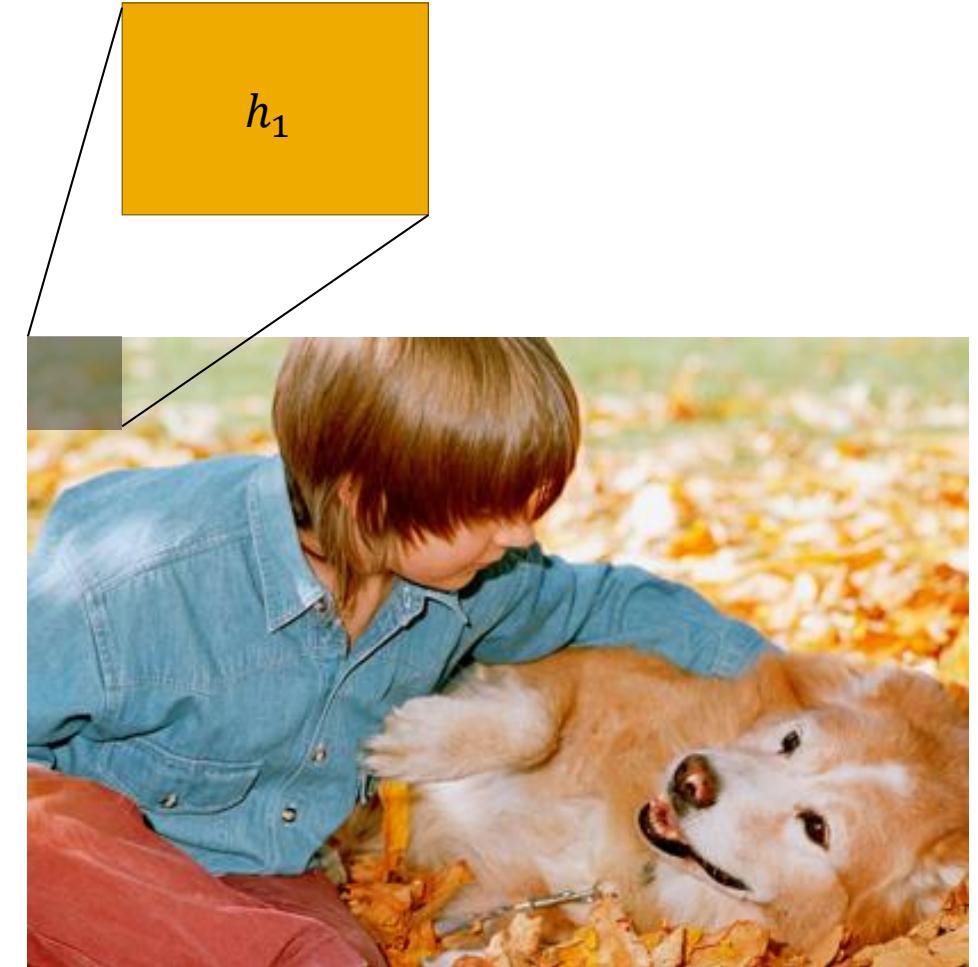
Each small patch is connected to the first hidden layer of the network

Contrast with a feed-forward neural network, which connects each pixel to each hidden node

Significantly more parameters:  
*image width × image height × size of hidden layer*

instead of:

*patch width × patch height × size of hidden layer*



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

In deeper layers of the network, the features become effectively larger and more abstract



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

In deeper layers of the network, the features become effectively larger and more abstract



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

In deeper layers of the network, the features become effectively larger and more abstract

## Advantages:

- Fewer parameters to tune
- Spatial invariance built in mathematically
- Retains spatial relationship between features



# Introduction to CNNs

Convolutional neural networks

CNNs apply *the same transformation* (feature maps) to each patch of the image

In deeper layers of the network, the features become effectively larger and more abstract

## Advantages:

- Fewer parameters to tune
- Spatial invariance built in mathematically
- Retains spatial relationship between features

The next lecture will explain CNNs and these advantages in detail

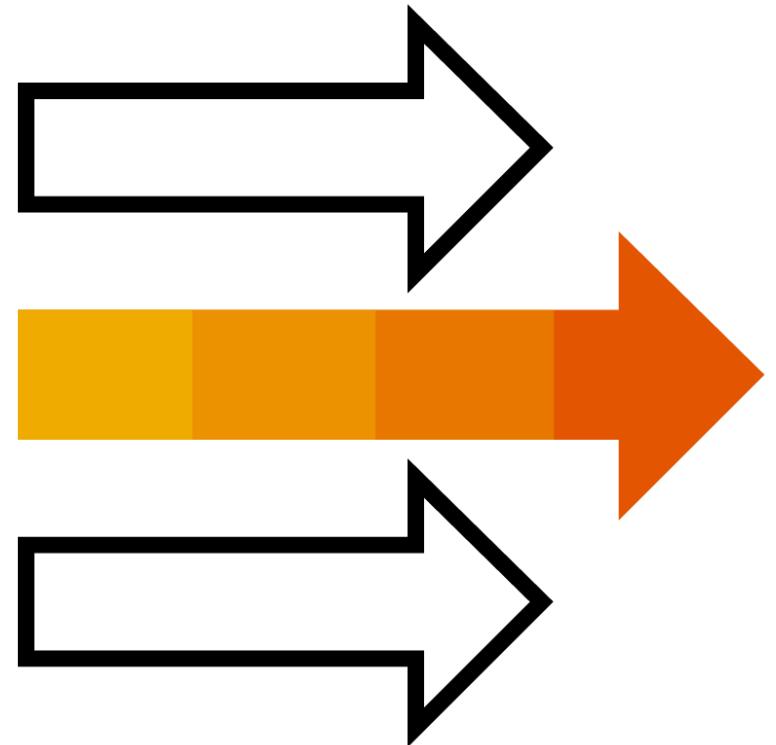


# Introduction to CNNs

Coming up next

## CNN Architecture I

- Convolutions
- Non-linearity
- Weight initialization



# Thank you.

Contact information:

[open@sap.com](mailto:open@sap.com)

# © 2017 SAP SE or an SAP affiliate company. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP SE or an SAP affiliate company.

The information contained herein may be changed without prior notice. Some software products marketed by SAP SE and its distributors contain proprietary software components of other software vendors. National product specifications may vary.

These materials are provided by SAP SE or an SAP affiliate company for informational purposes only, without representation or warranty of any kind, and SAP or its affiliated companies shall not be liable for errors or omissions with respect to the materials. The only warranties for SAP or SAP affiliate company products and services are those that are set forth in the express warranty statements accompanying such products and services, if any. Nothing herein should be construed as constituting an additional warranty.

In particular, SAP SE or its affiliated companies have no obligation to pursue any course of business outlined in this document or any related presentation, or to develop or release any functionality mentioned therein. This document, or any related presentation, and SAP SE's or its affiliated companies' strategy and possible future developments, products, and/or platform directions and functionality are all subject to change and may be changed by SAP SE or its affiliated companies at any time for any reason without notice. The information in this document is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. All forward-looking statements are subject to various risks and uncertainties that could cause actual results to differ materially from expectations. Readers are cautioned not to place undue reliance on these forward-looking statements, and they should not be relied upon in making purchasing decisions.

SAP and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP SE (or an SAP affiliate company) in Germany and other countries. All other product and service names mentioned are the trademarks of their respective companies.

See <http://global.sap.com/corporate-en/legal/copyright/index.epx> for additional trademark information and notices.



Week 4: Convolutional Networks

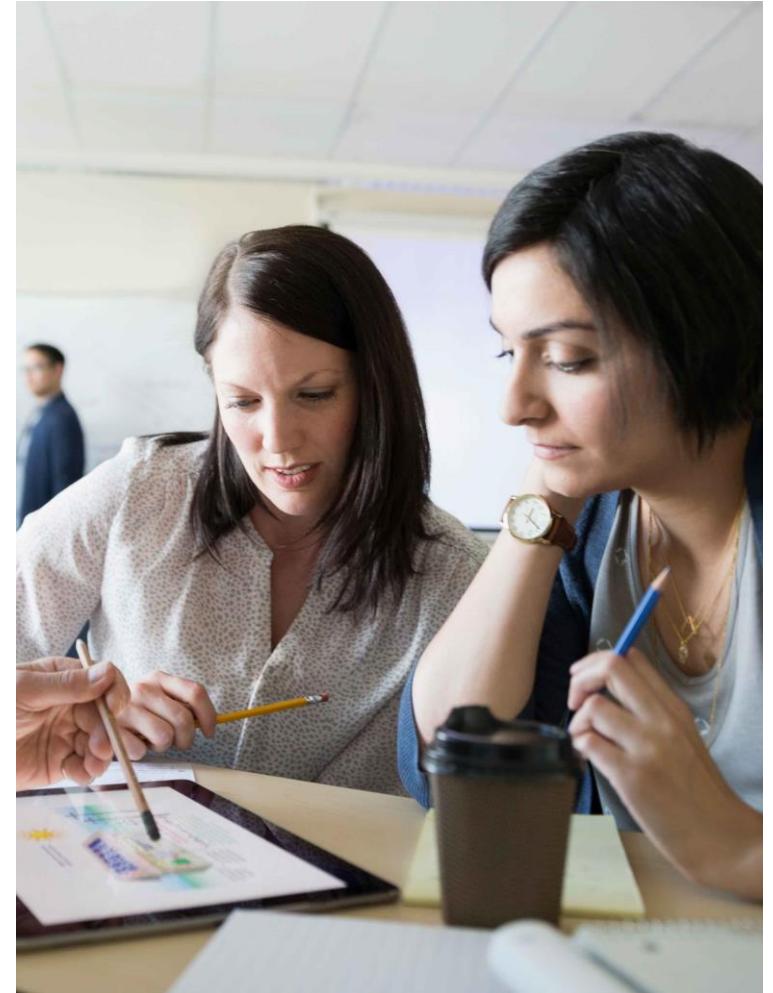
## **Unit 2: CNN Architecture Part I**

# CNN Architecture Part I

What we covered in the last unit

## Introduction to Convolutional Networks

- Biological inspiration
- Spatial invariance
- Basic idea: convolutional network

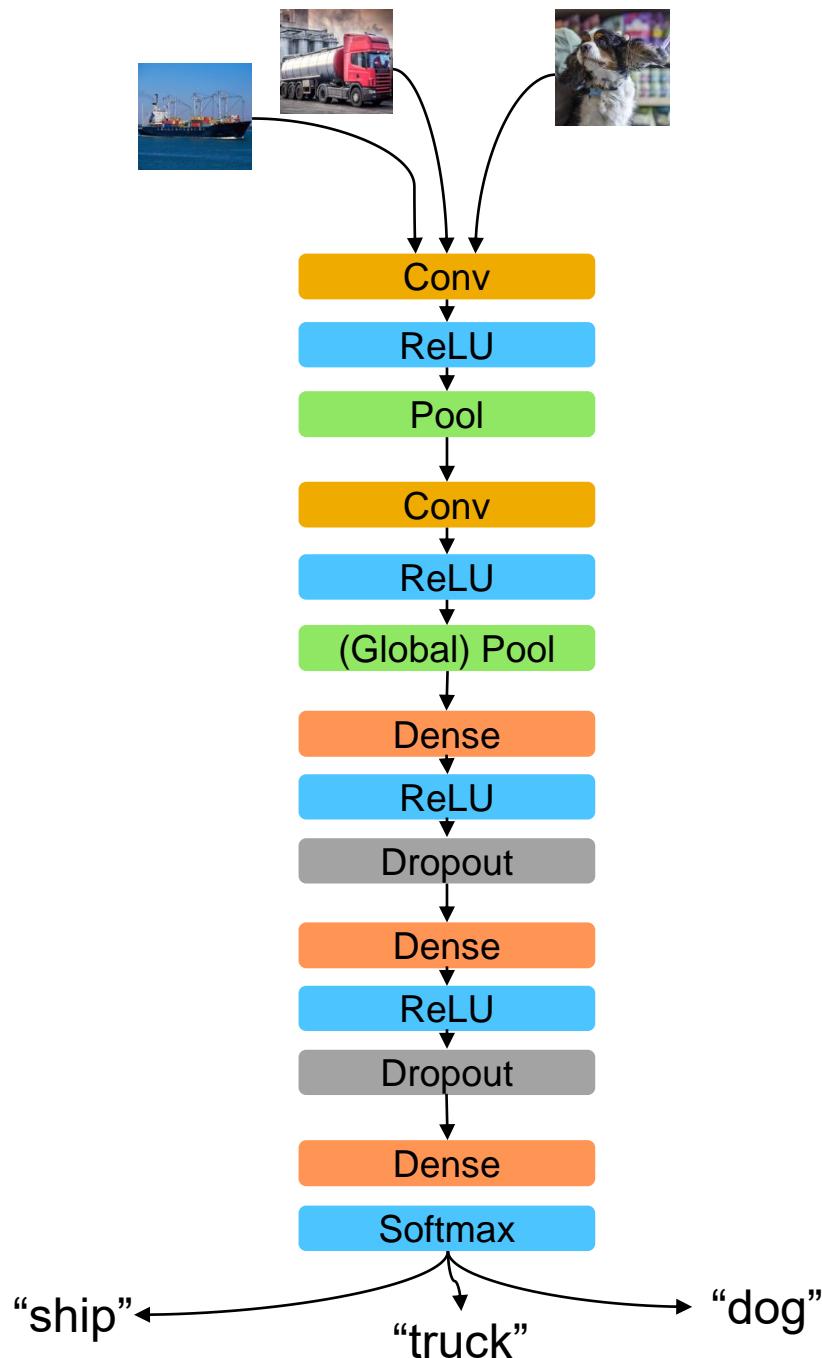


# CNN Architecture Part I & II

## Overview

### Content:

- Convolutions
- ReLU non-linearity (activation function)
- Weight initialization
- Pooling / global pooling
- Dense (fully connected) layers
- Dropout
- Softmax

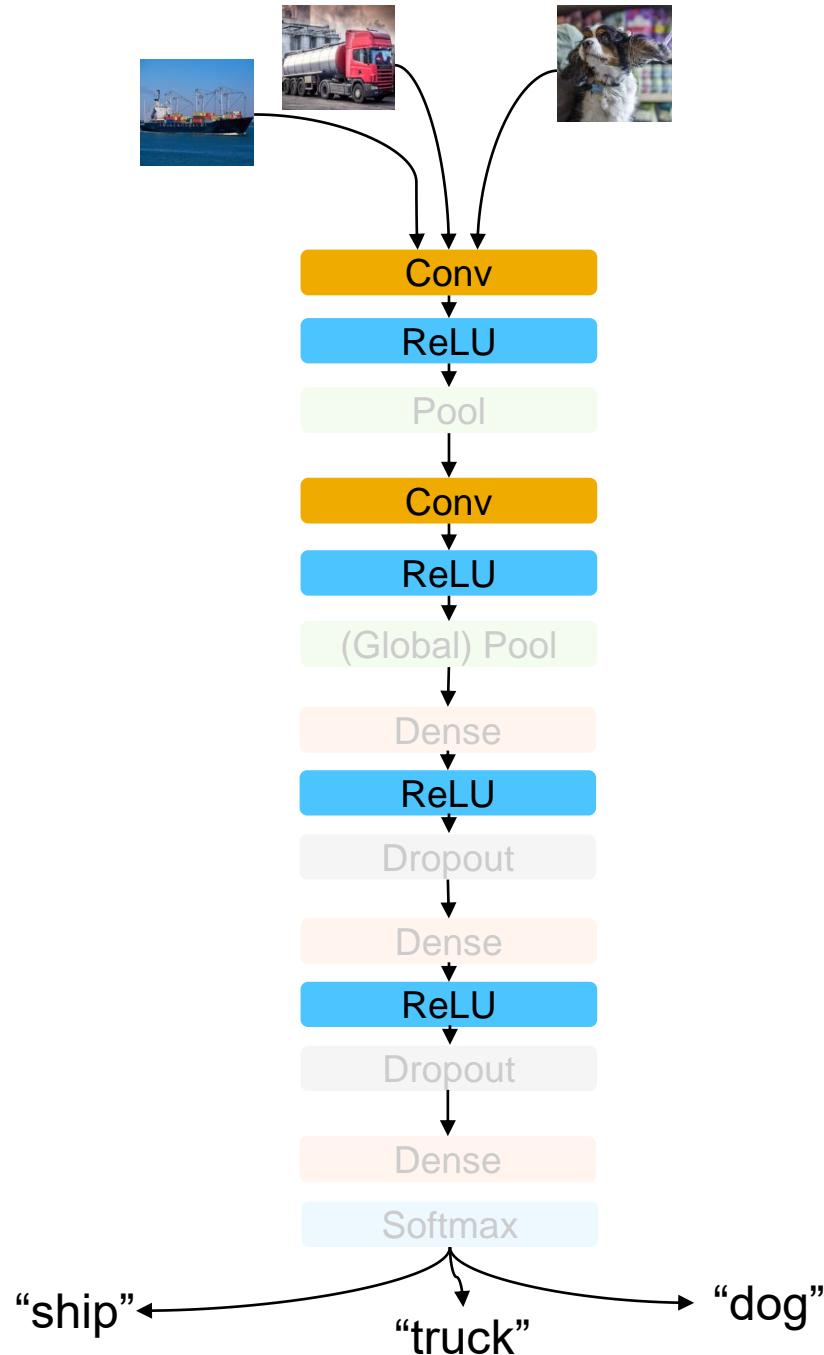


# CNN Architecture Part I

## Overview

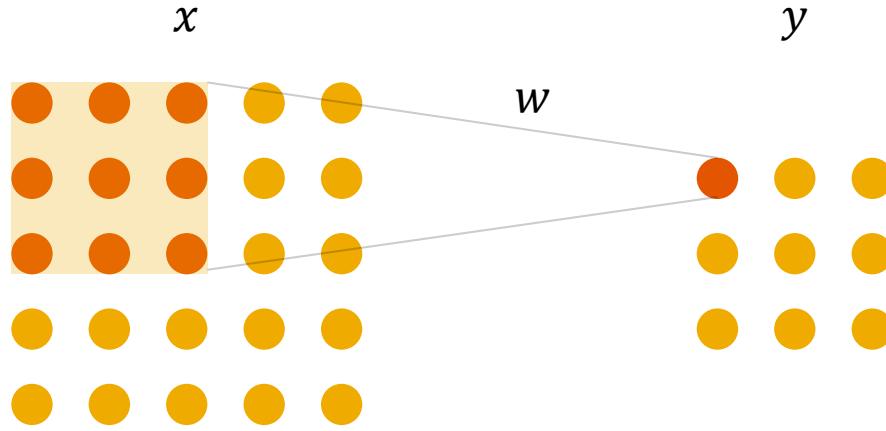
### Content:

- Convolutions
- ReLU non-linearity (activation function)
- Weight Initialization
- pooling / global pooling
- dense (fully connected) layers
- dropout
- softmax



# CNN Architecture Part I

## Convolutions



$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Input

Output

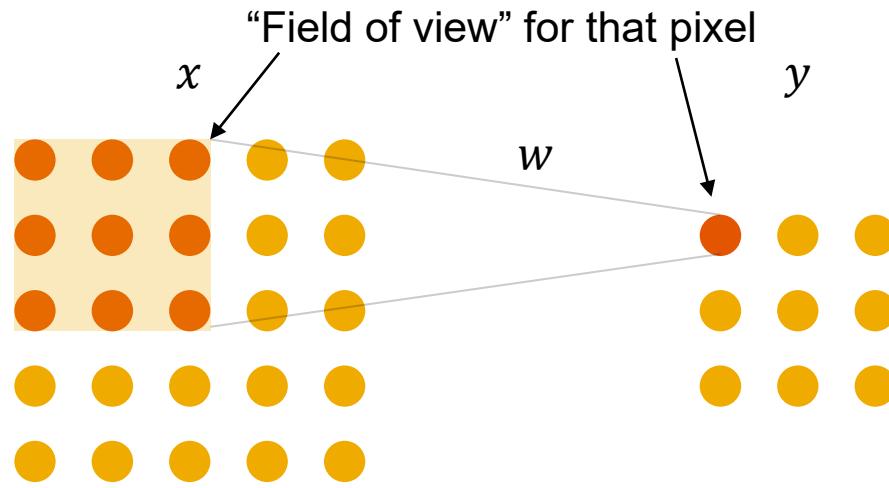
“valid” convolution

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



Input

Output

"valid" convolution

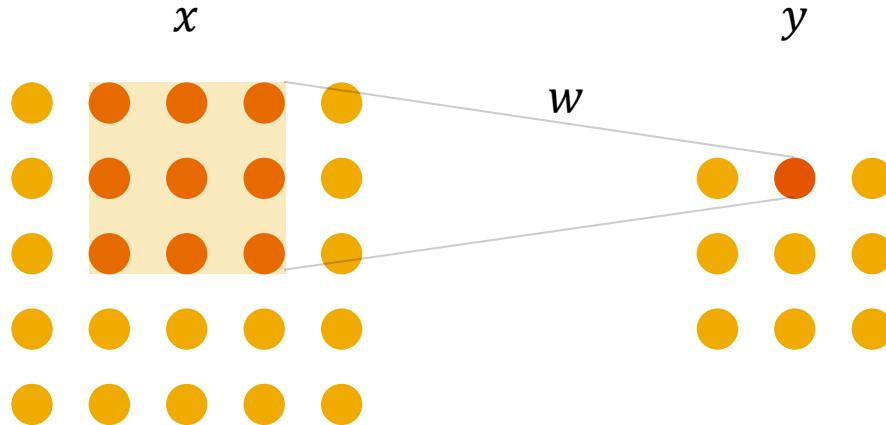
$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Input

Output

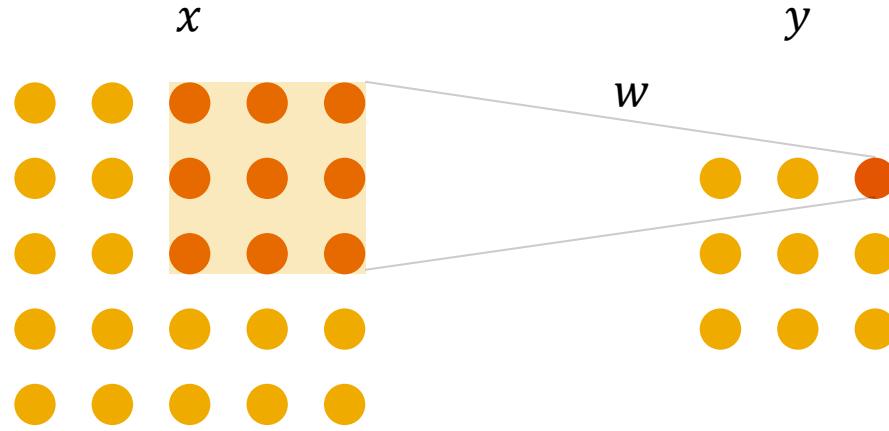
“valid” convolution

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Input

Output

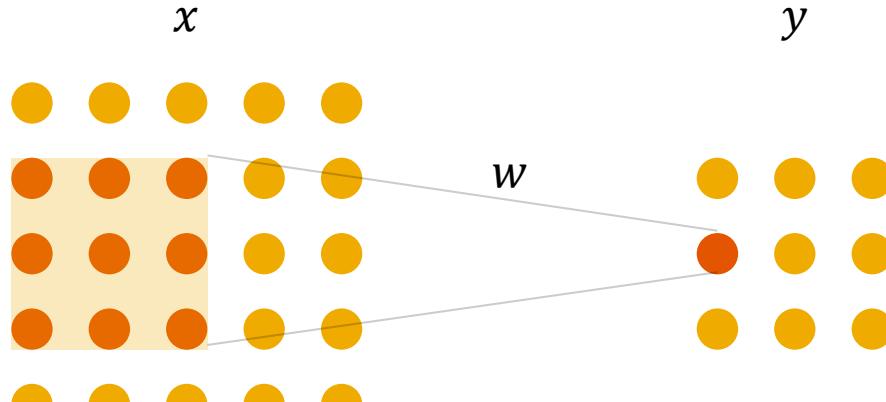
“valid” convolution

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Input

Output

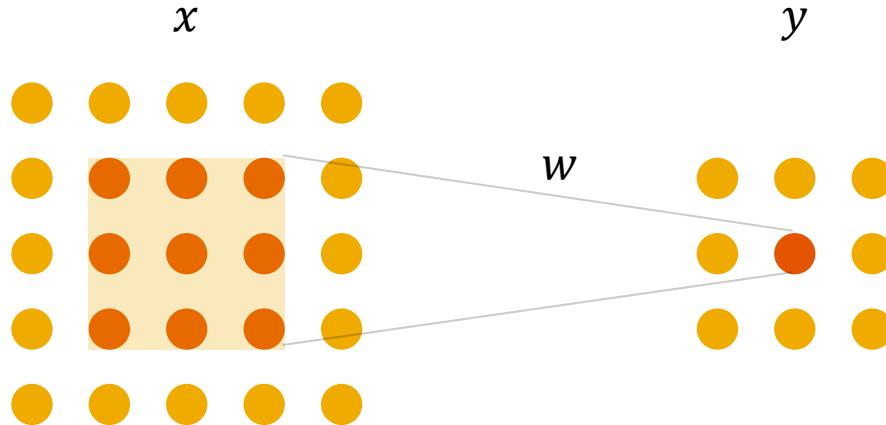
“valid” convolution

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



Input

Output

“valid” convolution

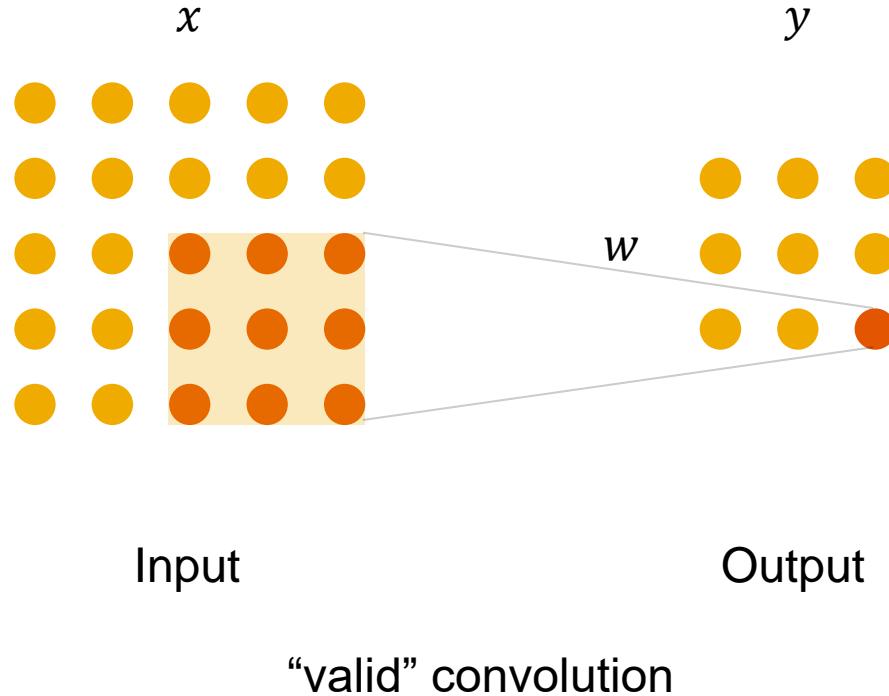
$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



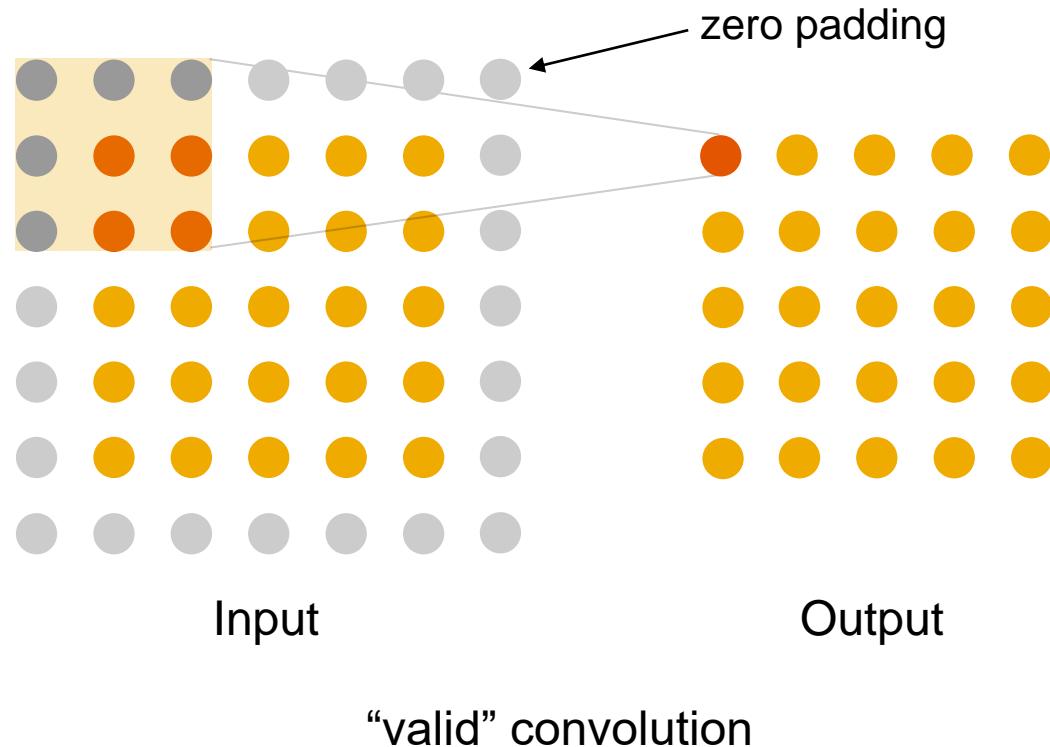
$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions



$$y_{ij} = \sum_{kl} w_{kl} x_{(k+i)(l+j)}$$

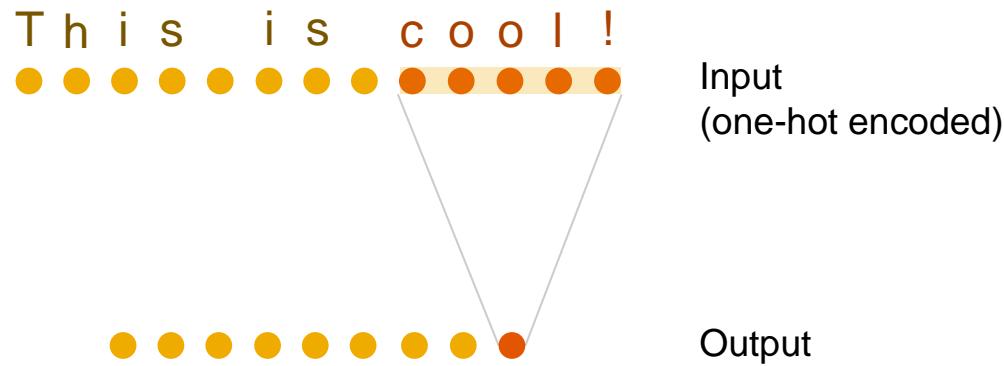
Weights are **shared** between pixels (neurons)!

- Imparts spatial invariance to network
- Reduces number of weights

# CNN Architecture Part I

## Convolutions in 1D

The same ideas generalize to 1D as well!



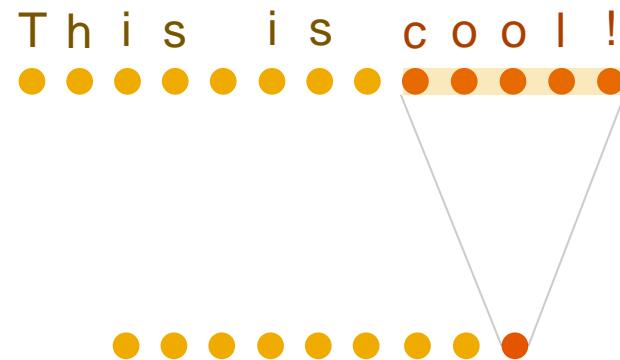
### 1D convolution:

e.g. audio signal

natural language processing

# CNN Architecture Part I

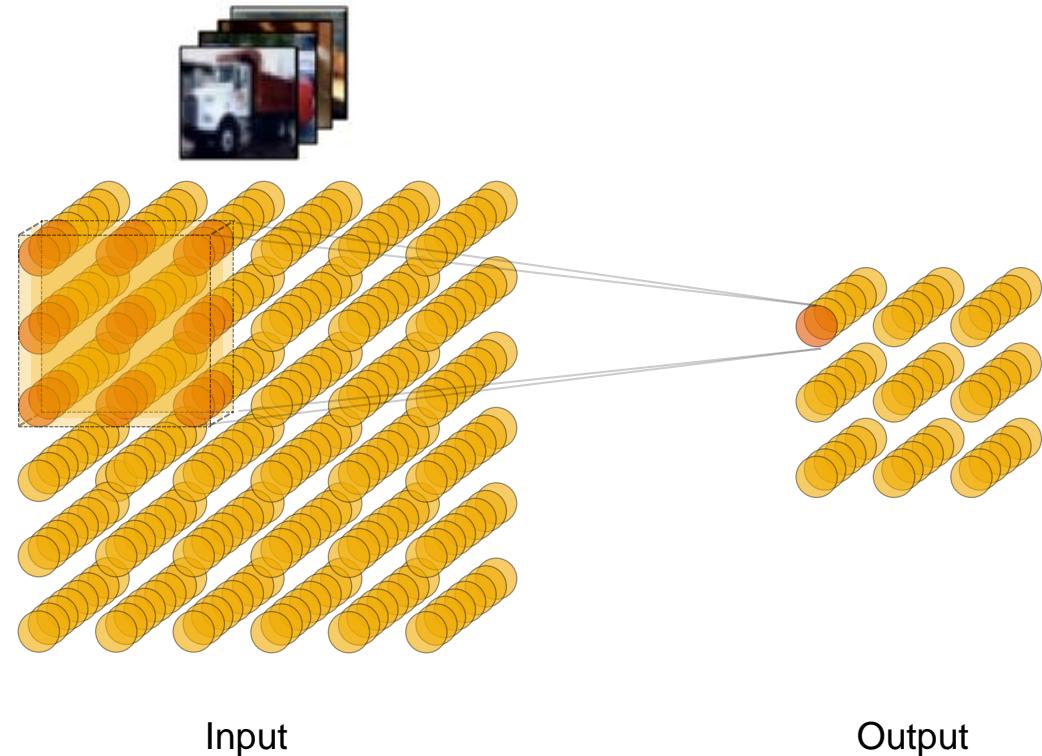
Convolutions in 1D & 3D



Input  
(one-hot encoded)

Output

The same ideas generalize to 1D & 3D as well!



Input

Output

## 1D convolution:

e.g. audio signal

natural language processing

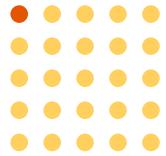
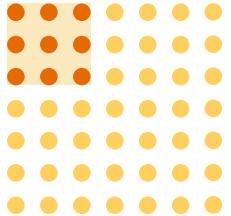
## 3D convolution:

e.g. video processing

# CNN Architecture Part I

Convolutions: popular kernel shapes

3x3 kernel

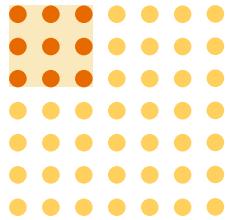


Most widely used in  
deep models

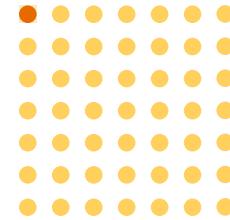
# CNN Architecture Part I

Convolutions: popular kernel shapes

3x3 kernel



1x1 kernel



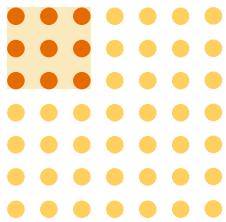
Most widely used in  
deep models

Will become  
clear on next slide

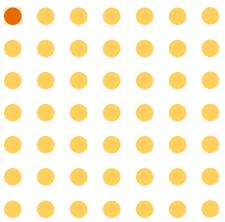
# CNN Architecture Part I

## Convolutions: popular kernel shapes

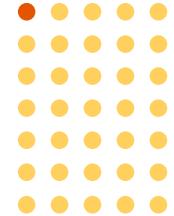
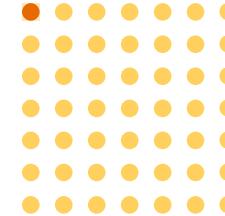
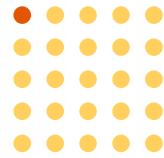
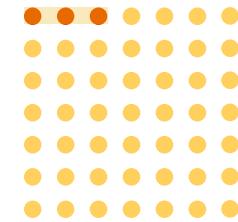
3x3 kernel



1x1 kernel



1x3 kernel  
(asymmetric)



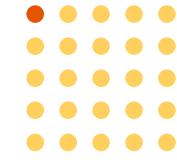
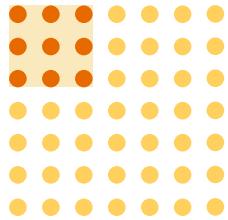
Most widely used in  
deep models

Will become  
clear on next slide

# CNN Architecture Part I

## Convolutions: popular kernel shapes

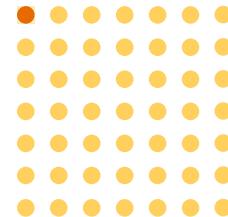
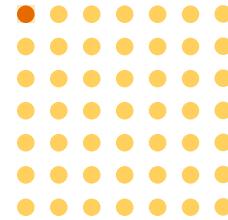
3x3 kernel



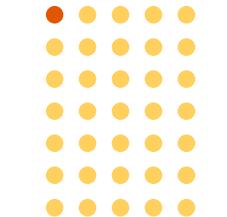
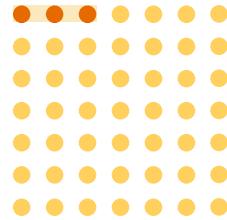
Most widely used in  
deep models

Will become  
clear on next slide

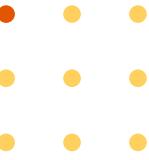
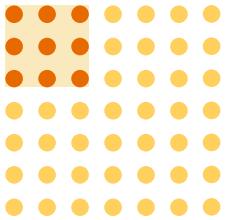
1x1 kernel



1x3 kernel  
(asymmetric)



3x3 kernel  
stride = 2

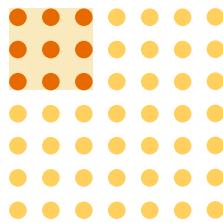


Downsampling

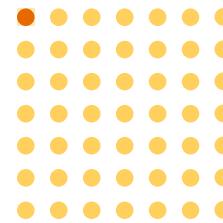
# CNN Architecture Part I

## Convolutions: popular kernel shapes

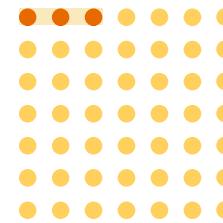
3x3 kernel



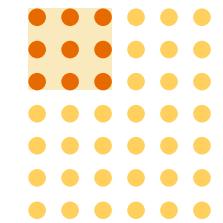
1x1 kernel



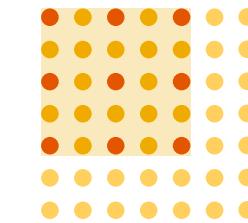
1x3 kernel  
(asymmetric)



3x3 kernel  
stride = 2



3x3 kernel  
dilation = 2



Most widely used in  
deep models

Will become  
clear on next slide

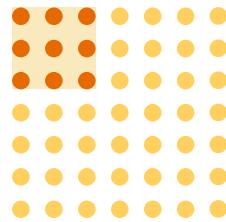
Downsampling

Increased field of view,  
but same complexity

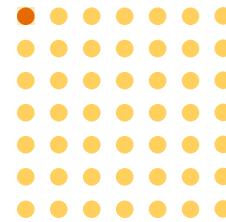
# CNN Architecture Part I

## Convolutions: popular kernel shapes

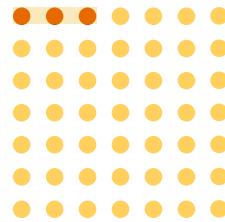
3x3 kernel



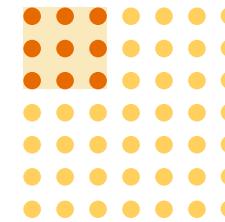
1x1 kernel



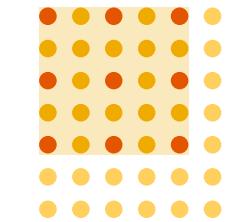
1x3 kernel  
(asymmetric)



3x3 kernel  
stride = 2



3x3 kernel  
dilation = 2



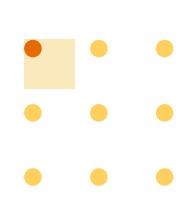
Most widely used in  
deep models

Will become  
clear on next slide

Downsampling

Increased field of view,  
but same complexity

2x2 kernel  
stride = 2



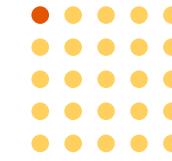
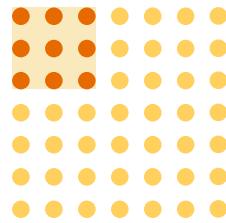
Upsampling

Transposed Convolution  
aka deconvolution  
aka fractionally strided convolution

# CNN Architecture Part I

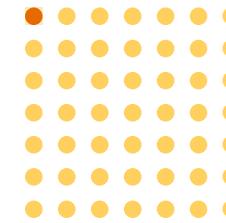
## Convolutions: popular kernel shapes

3x3 kernel



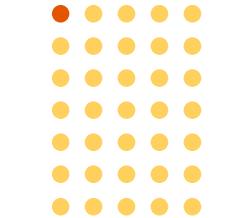
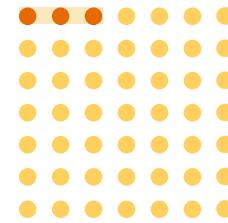
Most widely used in  
deep models

1x1 kernel

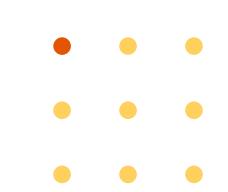
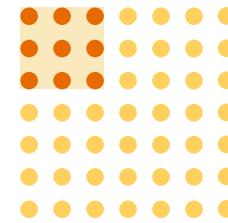


Will become  
clear on next slide

1x3 kernel  
(asymmetric)

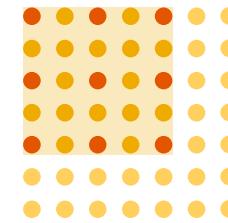


3x3 kernel  
stride = 2



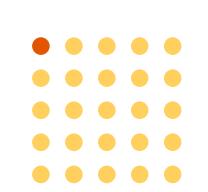
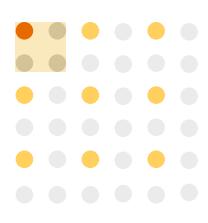
Downsampling

3x3 kernel  
dilation = 2



Increased field of view,  
but same complexity

2x2 kernel  
stride = 2

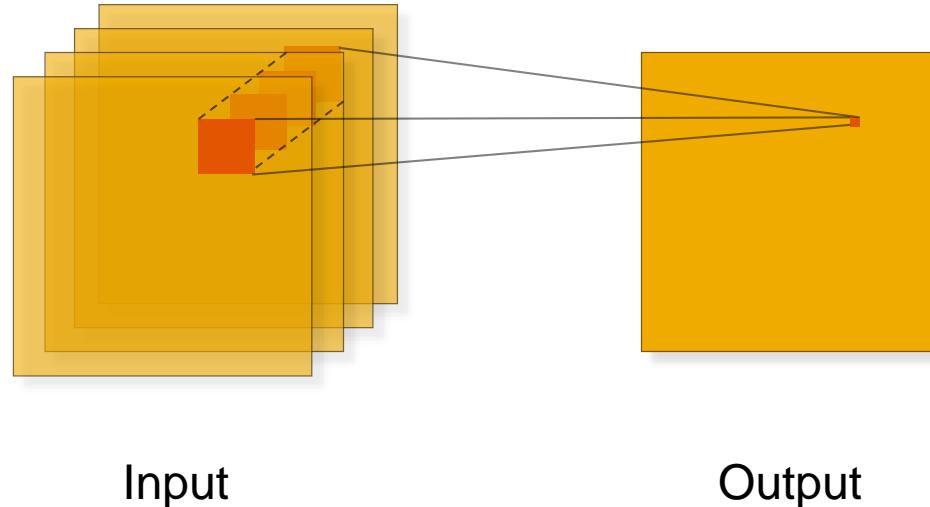


Upsampling

Transposed Convolution  
aka deconvolution  
aka fractionally strided convolution

# CNN Architecture Part I

## Convolutions

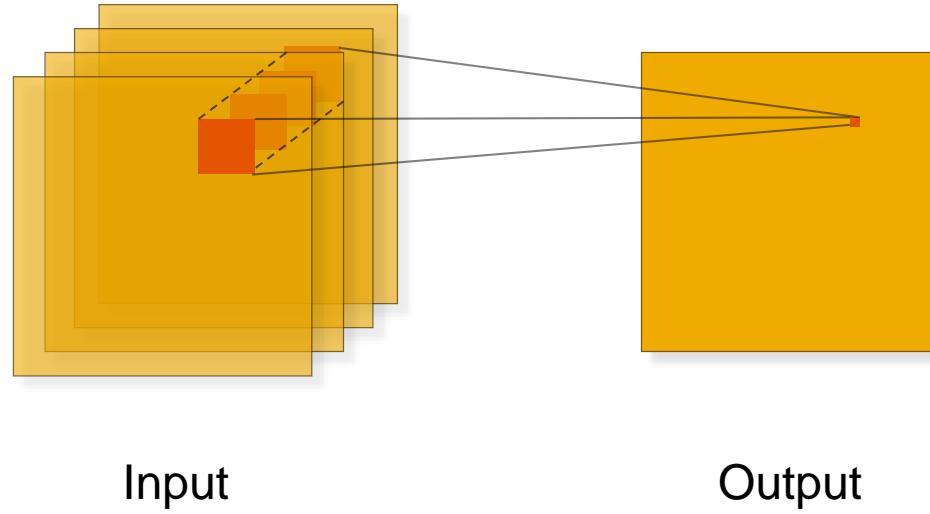


$$y_{ij} = \sum_C \sum_{kl} w_{Ckl} x_{C(k+i)(l+j)}$$

We are not just summing over pixels, but also across input channels!

# CNN Architecture Part I

## Convolutions



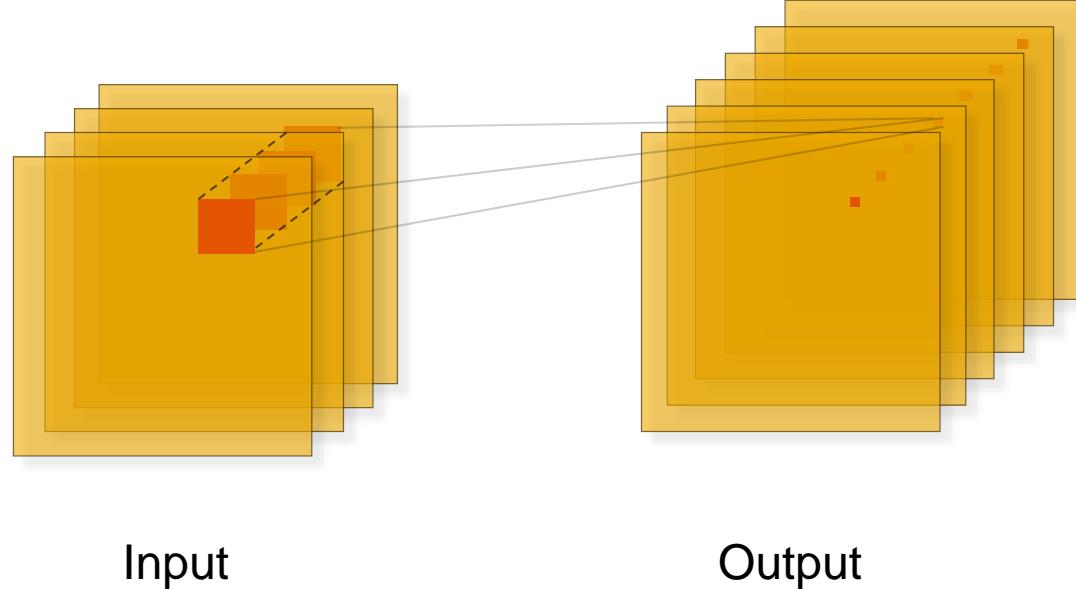
$$y_{ij} = \sum_C \sum_{kl} w_{Ckl} x_{C(k+i)(l+j)}$$

One “shared” set of weights  
per input channel

We are not just summing over pixels, but also across input channels!

# CNN Architecture Part I

## Convolutions



Multiple output channels

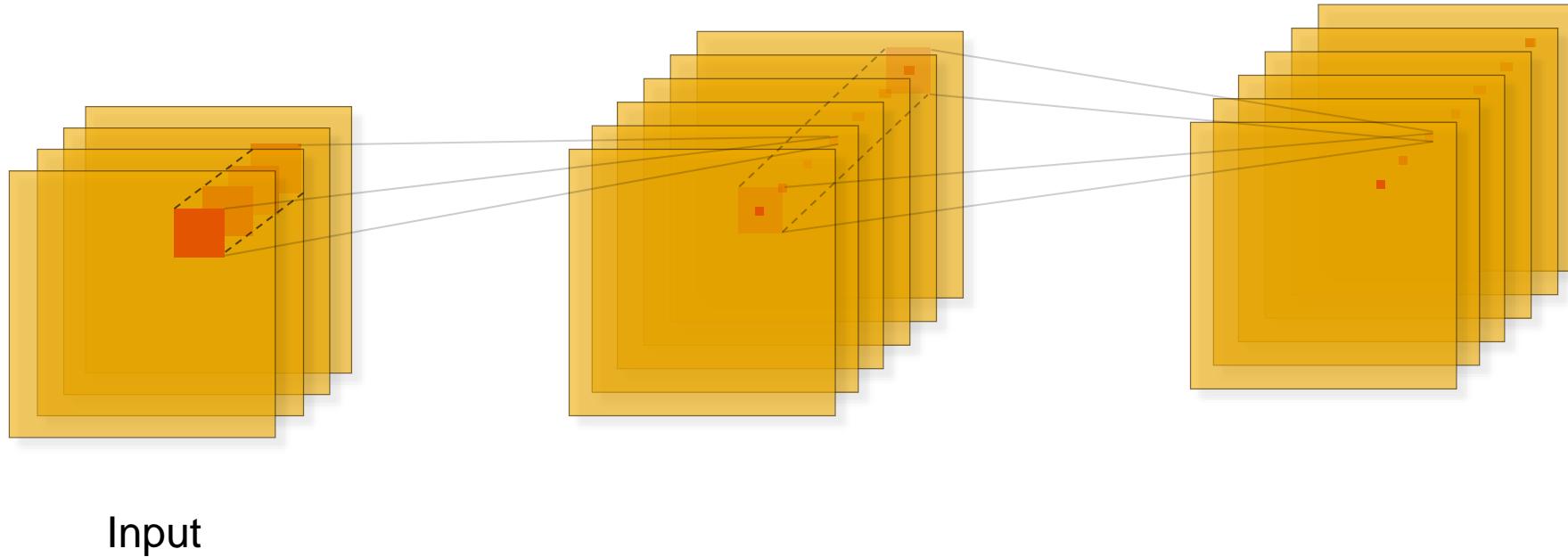
$$y_{ij}^I = \sum_C \sum_{kl} w_{Ckl}^I x_{C(k+i)(l+j)}$$

One "shared" set of  
weights per input channel  
and per output channel

We can have multiple output channels (feature maps) as well!

# CNN Architecture Part I

Convolutions: going deeper

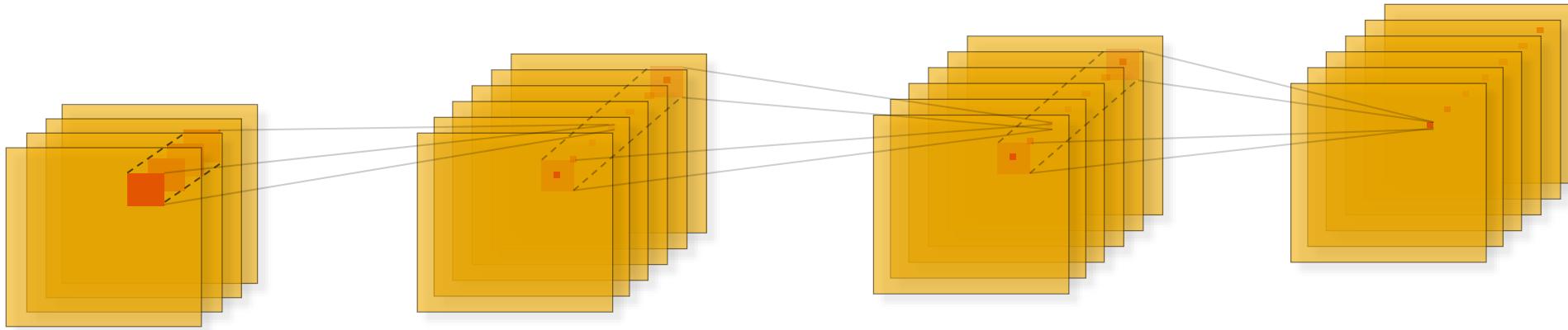


Input

We can chain multiple convolutions to build a deep convolutional neural network!

# CNN Architecture Part I

Convolutions: going deeper

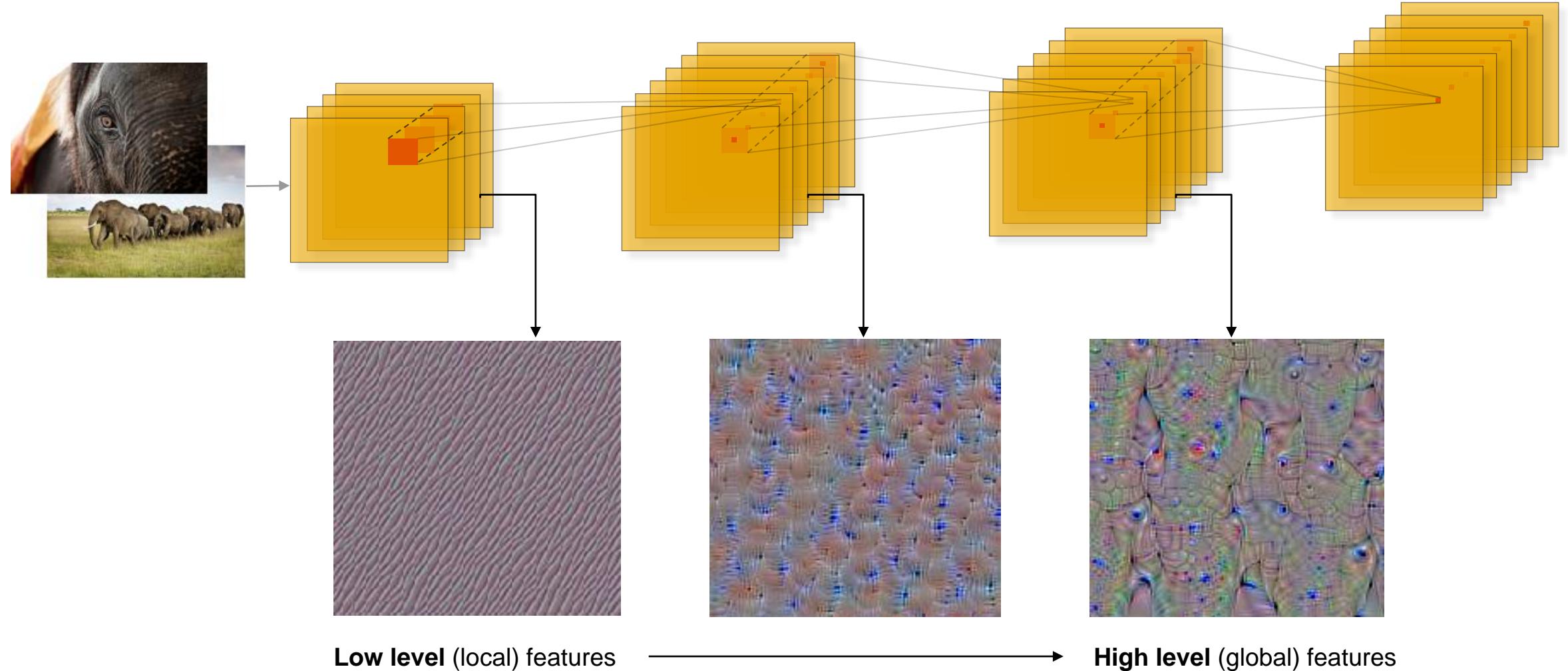


Input

We can chain multiple convolutions to build a deep convolutional neural network!

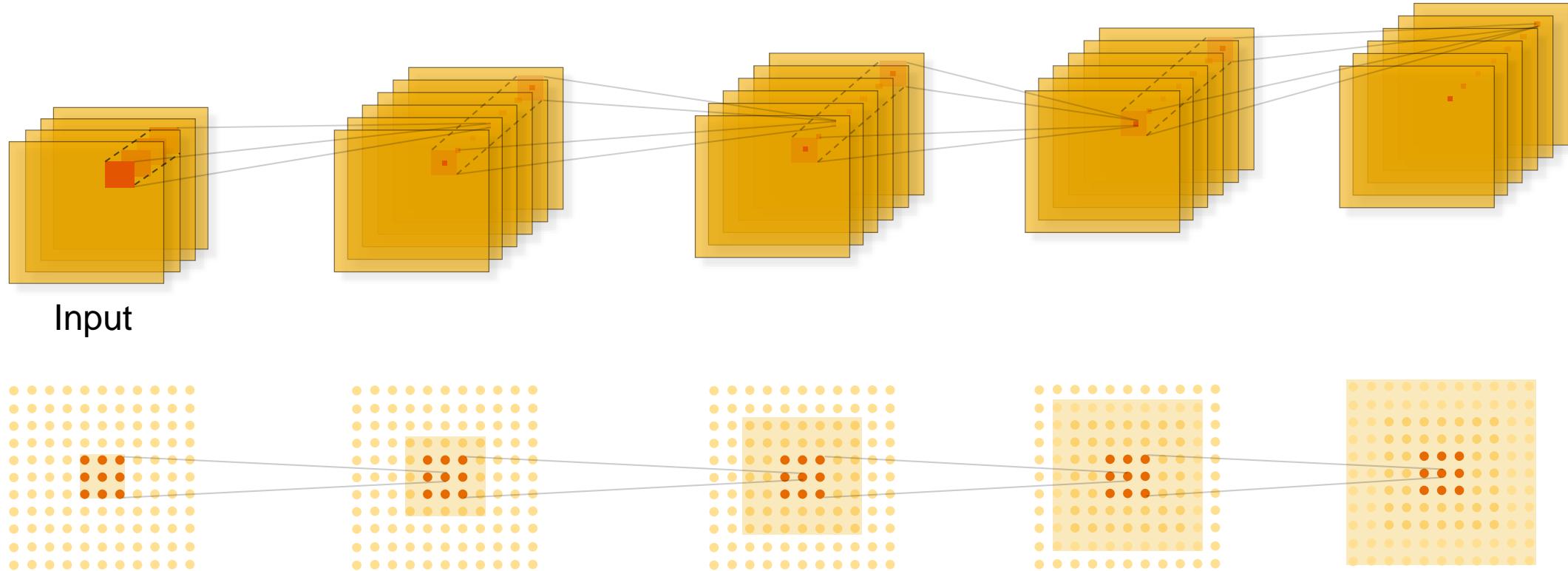
# CNN Architecture Part I

Convolutions: going deeper



# CNN Architecture Part I

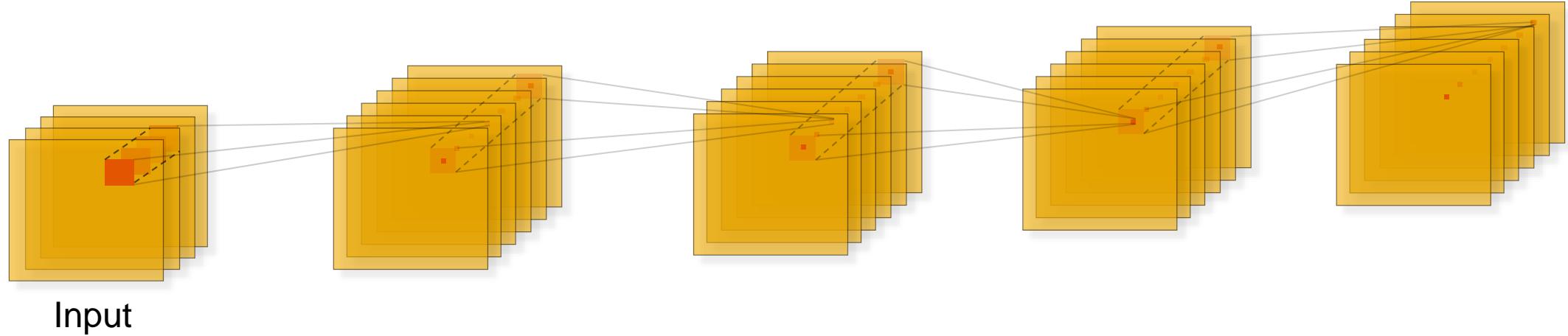
Convolutions: going deeper



Note: After each layer, the effective field of view *within the input image* is increased!

# CNN Architecture Part I

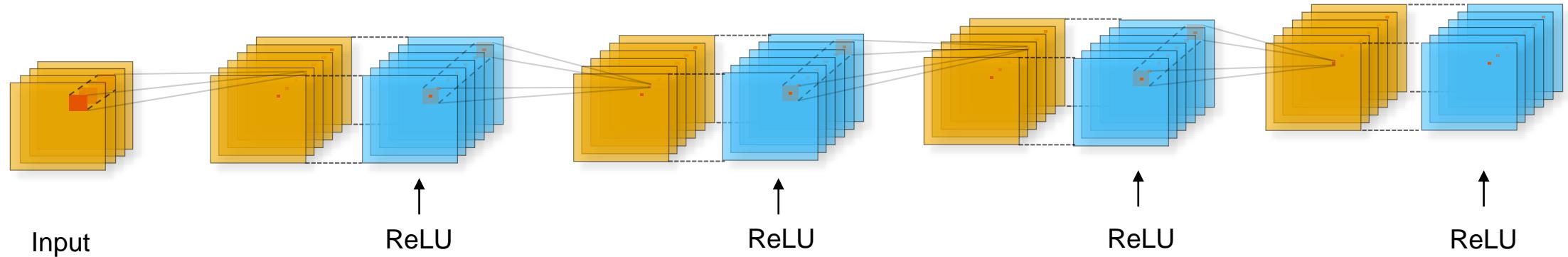
Non-linearity: ReLUs



After each convolution, we want to have a non-linearity to increase expressive power!

# CNN Architecture Part I

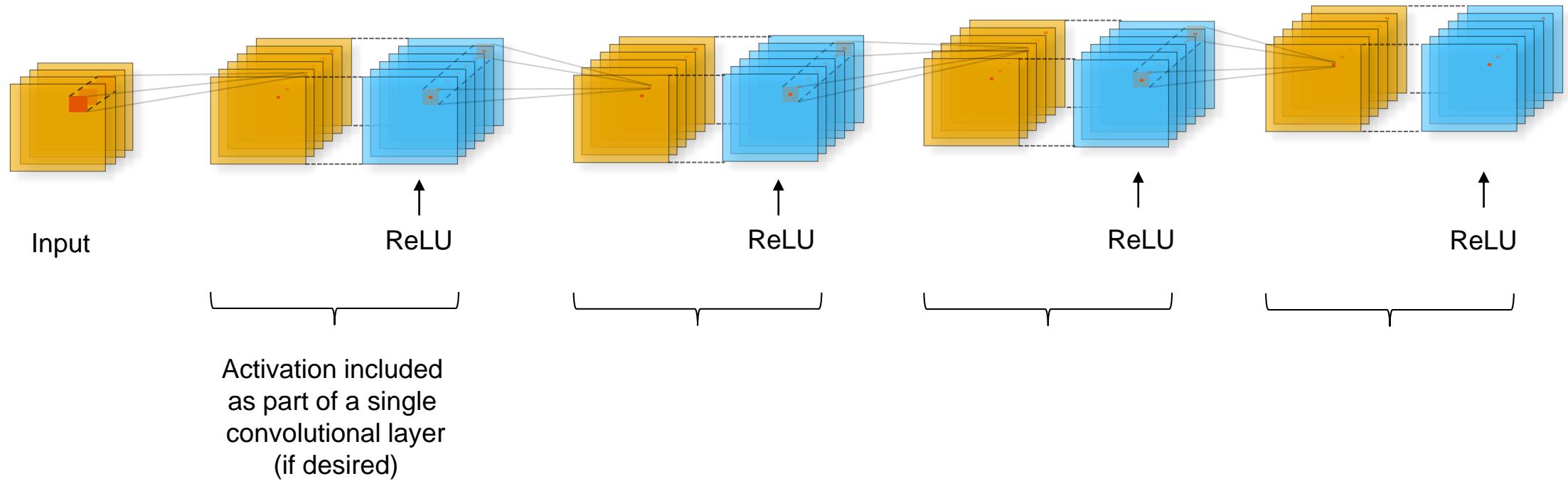
Non-linearity: ReLUs



After each convolution, we want to have a **non-linearity** to increase expressive power!

# CNN Architecture Part I

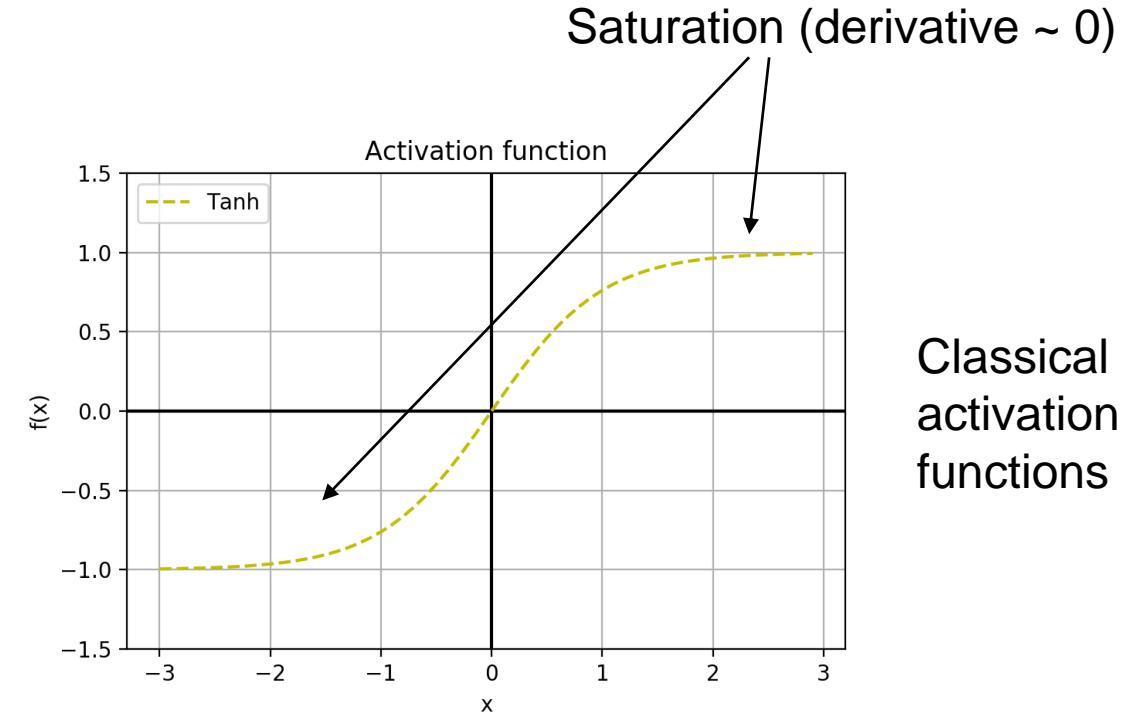
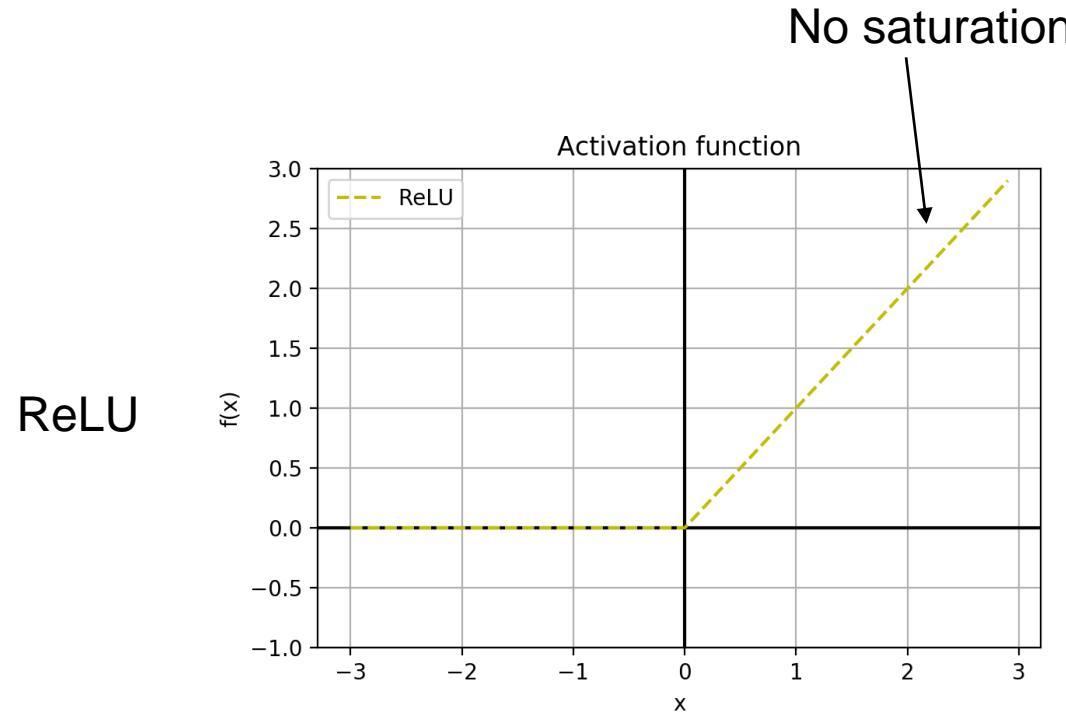
Non-linearity: ReLUs



After each convolution, we want to have a **non-linearity** to increase expressive power!

# CNN Architecture Part I

Non-linearity: rectified linear unit (ReLU)



Classical activation functions

- ReLUs do not **saturate**

(reduced vanishing gradient problem; see RNNs, Week 3)

- Induce sparsity

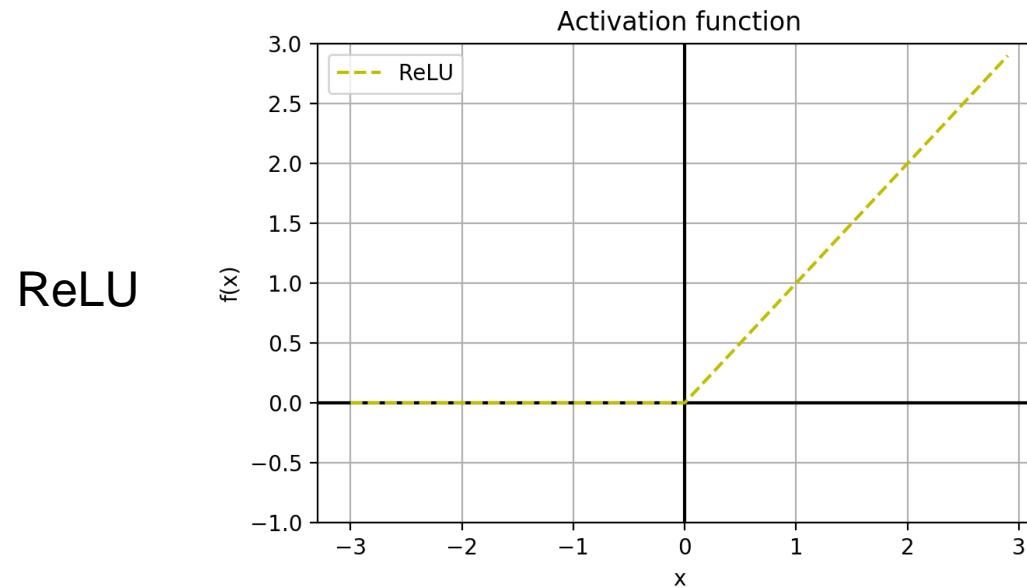
Deep Sparse Rectifier Neural Networks

Xavier Glorot, Antoine Bordes, Yoshua Bengio;

*Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, PMLR 15:315-323, 2011*

# CNN Architecture Part I

Non-linearity: parametric rectified linear unit (PReLU)

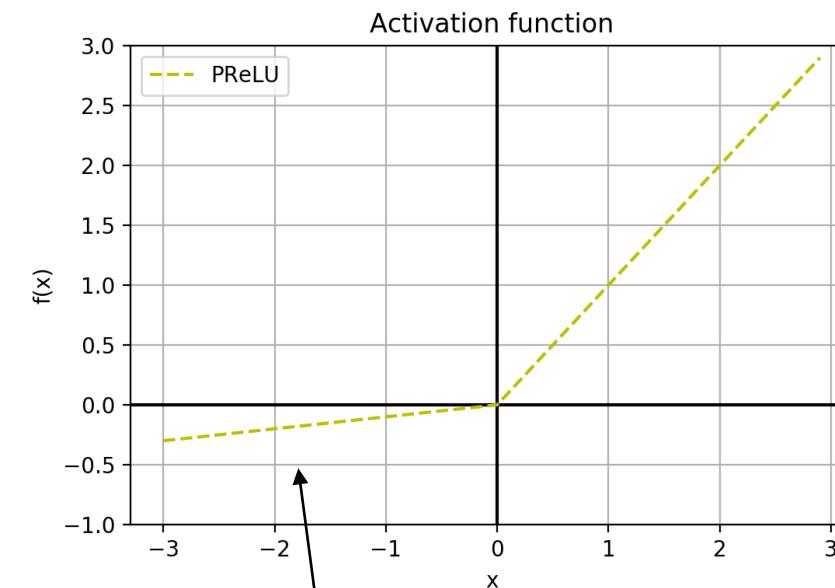


- ReLUs do not saturate
- Induce sparsity
- But they may “die” during training (never ever activate again)

**Delving Deep into Rectifiers:  
Surpassing Human-Level Performance on ImageNet Classification**

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

ProceedingICCV '15 Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV) Pages 1026-1034



Parametric (leaky) ReLU (PReLU)  
does not “die”

# CNN Architecture Part I

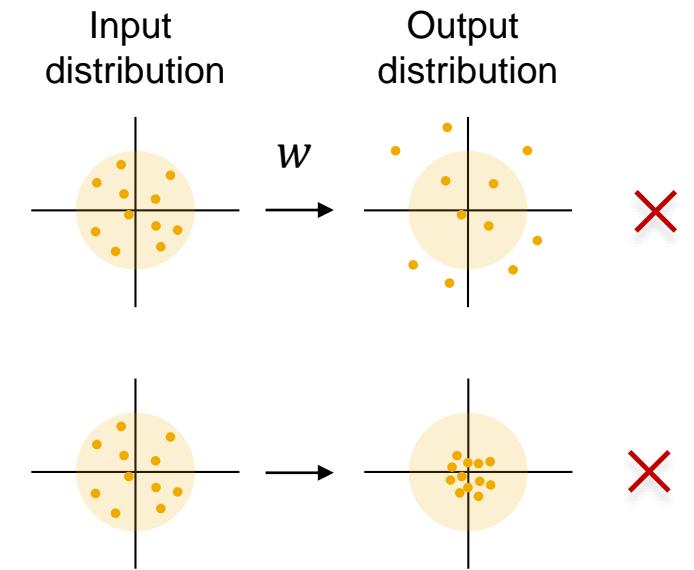
## Weight initialization

- Weights in each layer are **randomly initialized**

# CNN Architecture Part I

## Weight initialization

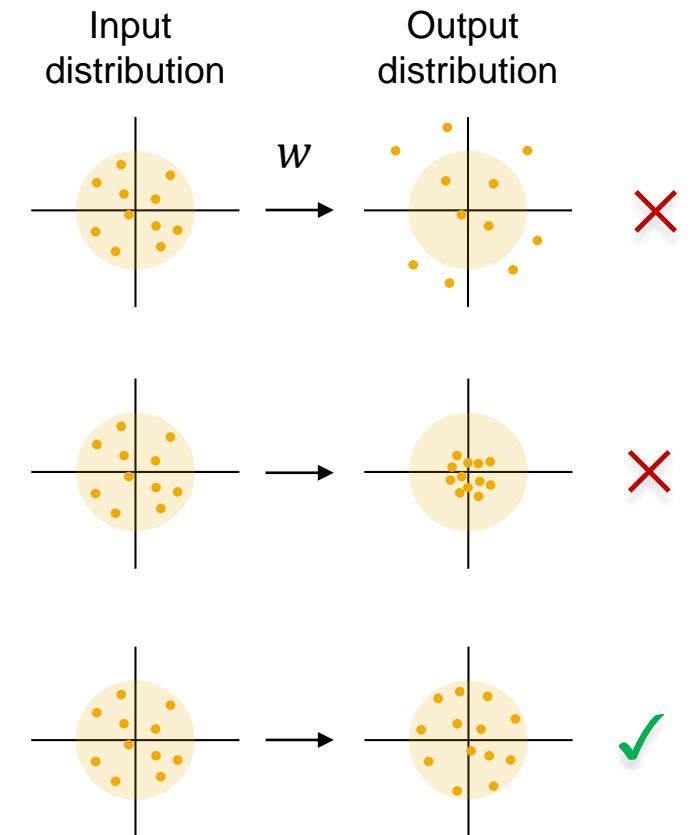
- Weights in each layer are **randomly initialized**
- But they should be initialized such that an input signal is ***neither amplified nor damped!***



# CNN Architecture Part I

## Weight initialization

- Weights in each layer are **randomly initialized**
- But they should be initialized such that an input signal is ***neither amplified nor damped!***
- This property **depends on the number of input/output connections, and the non-linearity**



# CNN Architecture Part I

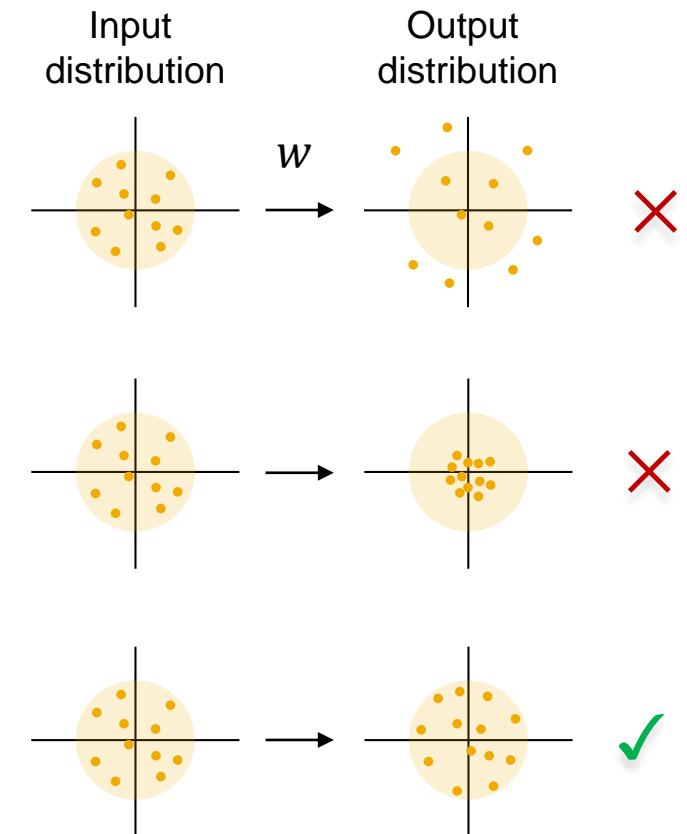
## Weight initialization

- Weights in each layer are **randomly initialized**
- But they should be initialized such that an input signal is **neither amplified nor damped!**
- This property **depends on the number of input/output connections, and the non-linearity**

For ReLU, use “**He’s**” method (aka Gaussian / variance scaling initialization)

$$w \sim P_{\text{Gaussian}}, \quad Var = \sqrt{\frac{2}{N}},$$

N: number of input or output channels, or average of the two



Delving Deep into Rectifiers:  
Surpassing Human-Level Performance on ImageNet Classification

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun

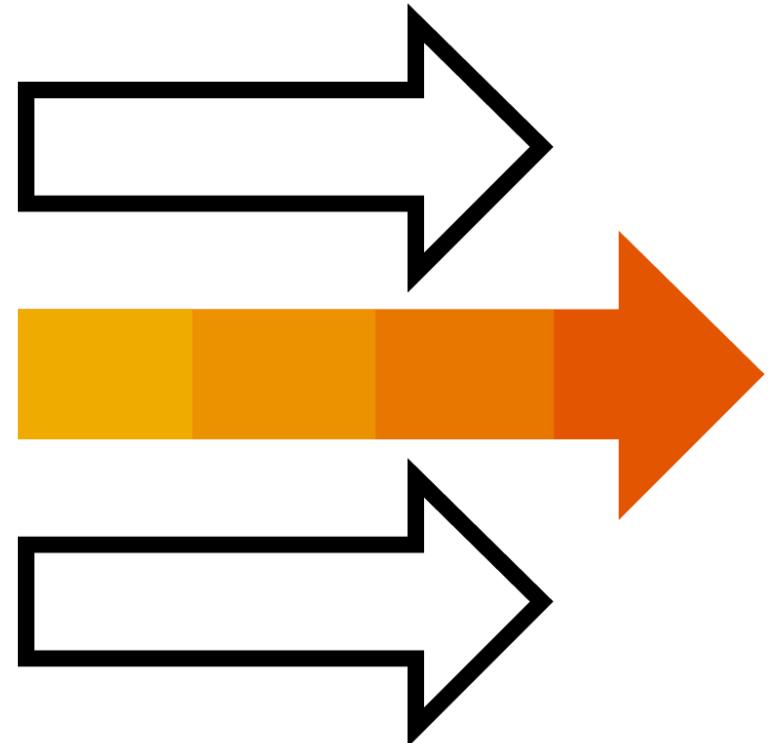
ProceedingICCV '15 Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV) Pages 1026-1034

# CNN Architecture Part I

Coming up next

## CNN Architecture II

- Pooling
- Dense layers
- Dropout
- Softmax



# Thank you.

Contact information:

[open@sap.com](mailto:open@sap.com)

# © 2017 SAP SE or an SAP affiliate company. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP SE or an SAP affiliate company.

The information contained herein may be changed without prior notice. Some software products marketed by SAP SE and its distributors contain proprietary software components of other software vendors. National product specifications may vary.

These materials are provided by SAP SE or an SAP affiliate company for informational purposes only, without representation or warranty of any kind, and SAP or its affiliated companies shall not be liable for errors or omissions with respect to the materials. The only warranties for SAP or SAP affiliate company products and services are those that are set forth in the express warranty statements accompanying such products and services, if any. Nothing herein should be construed as constituting an additional warranty.

In particular, SAP SE or its affiliated companies have no obligation to pursue any course of business outlined in this document or any related presentation, or to develop or release any functionality mentioned therein. This document, or any related presentation, and SAP SE's or its affiliated companies' strategy and possible future developments, products, and/or platform directions and functionality are all subject to change and may be changed by SAP SE or its affiliated companies at any time for any reason without notice. The information in this document is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. All forward-looking statements are subject to various risks and uncertainties that could cause actual results to differ materially from expectations. Readers are cautioned not to place undue reliance on these forward-looking statements, and they should not be relied upon in making purchasing decisions.

SAP and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP SE (or an SAP affiliate company) in Germany and other countries. All other product and service names mentioned are the trademarks of their respective companies.

See <http://global.sap.com/corporate-en/legal/copyright/index.epx> for additional trademark information and notices.



Week 4: Convolutional Networks

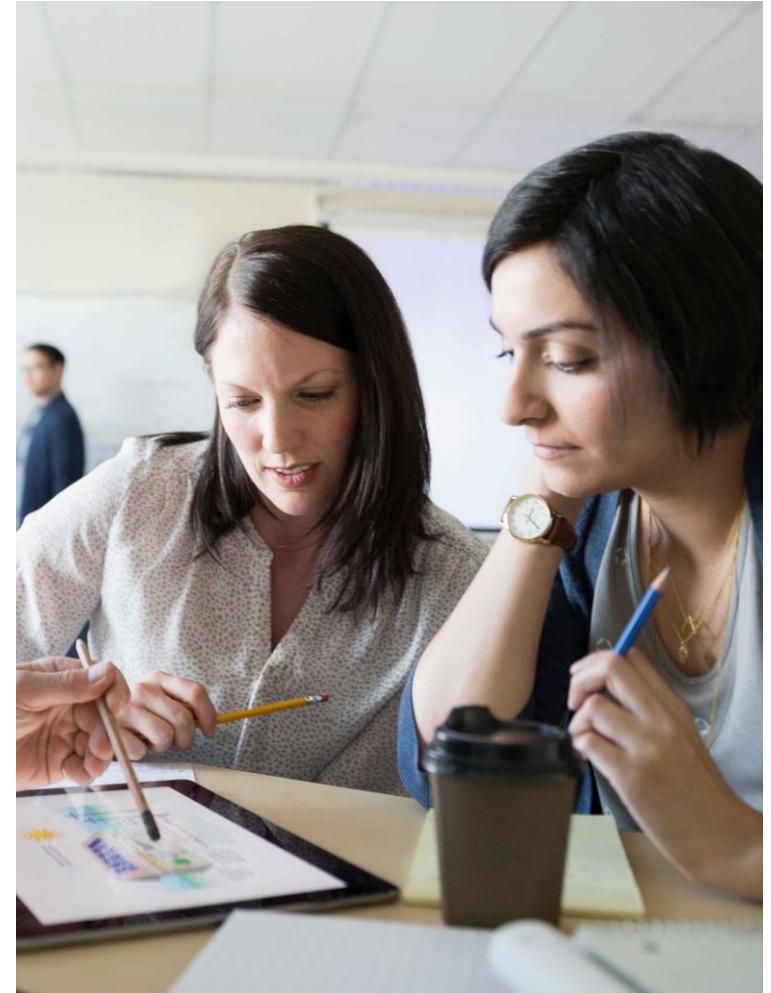
## **Unit 3: CNN Architecture Part II**

# CNN Architecture Part II

What we covered in the last unit

## CNN Architecture Part I

- Convolutions
- Non-linearity
- Weight initialization

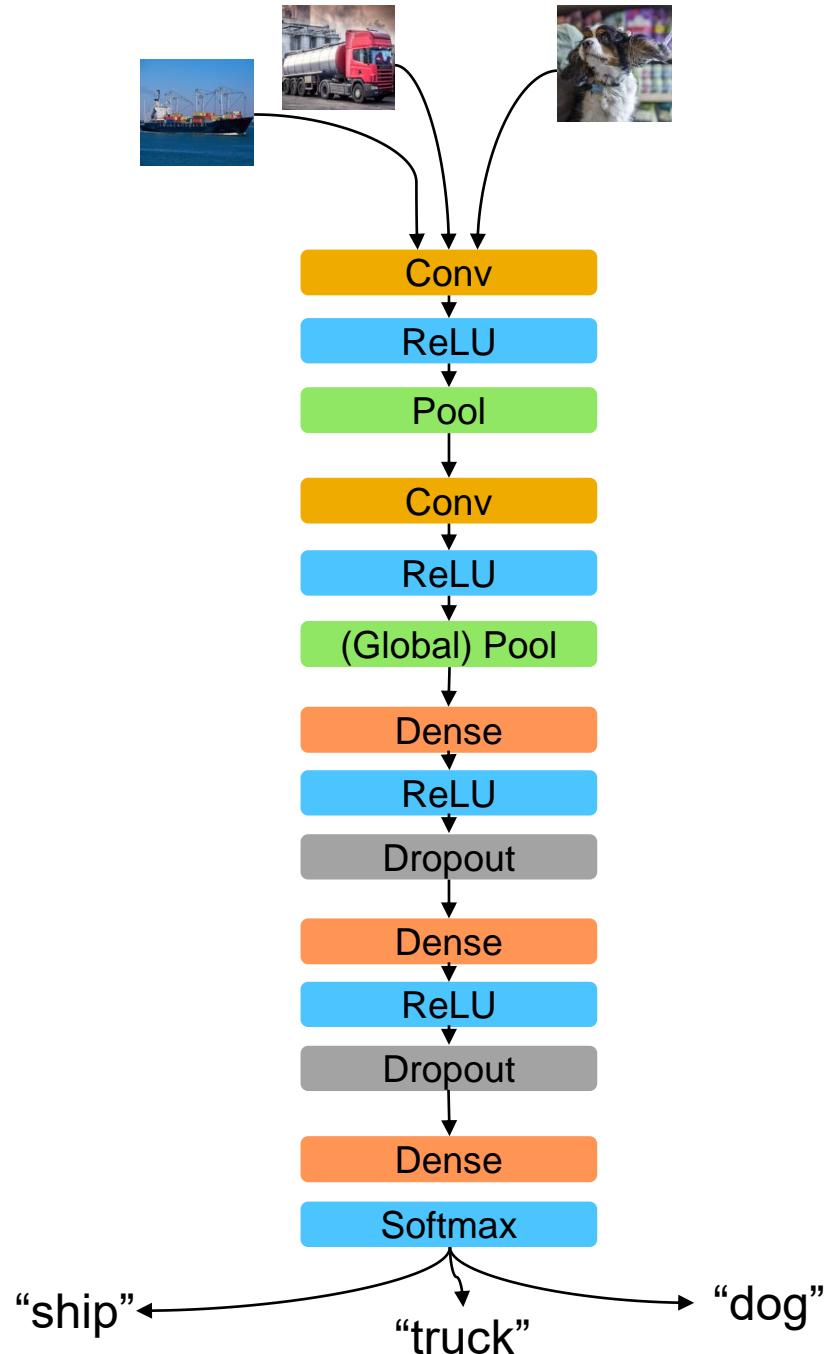


# CNN Architecture Part I & II

## Overview

### Content:

- Convolutions
- ReLU non-linearity (activation function)
- Weight initialization
- Pooling / global pooling
- Dense (fully connected) layers
- Dropout
- Softmax

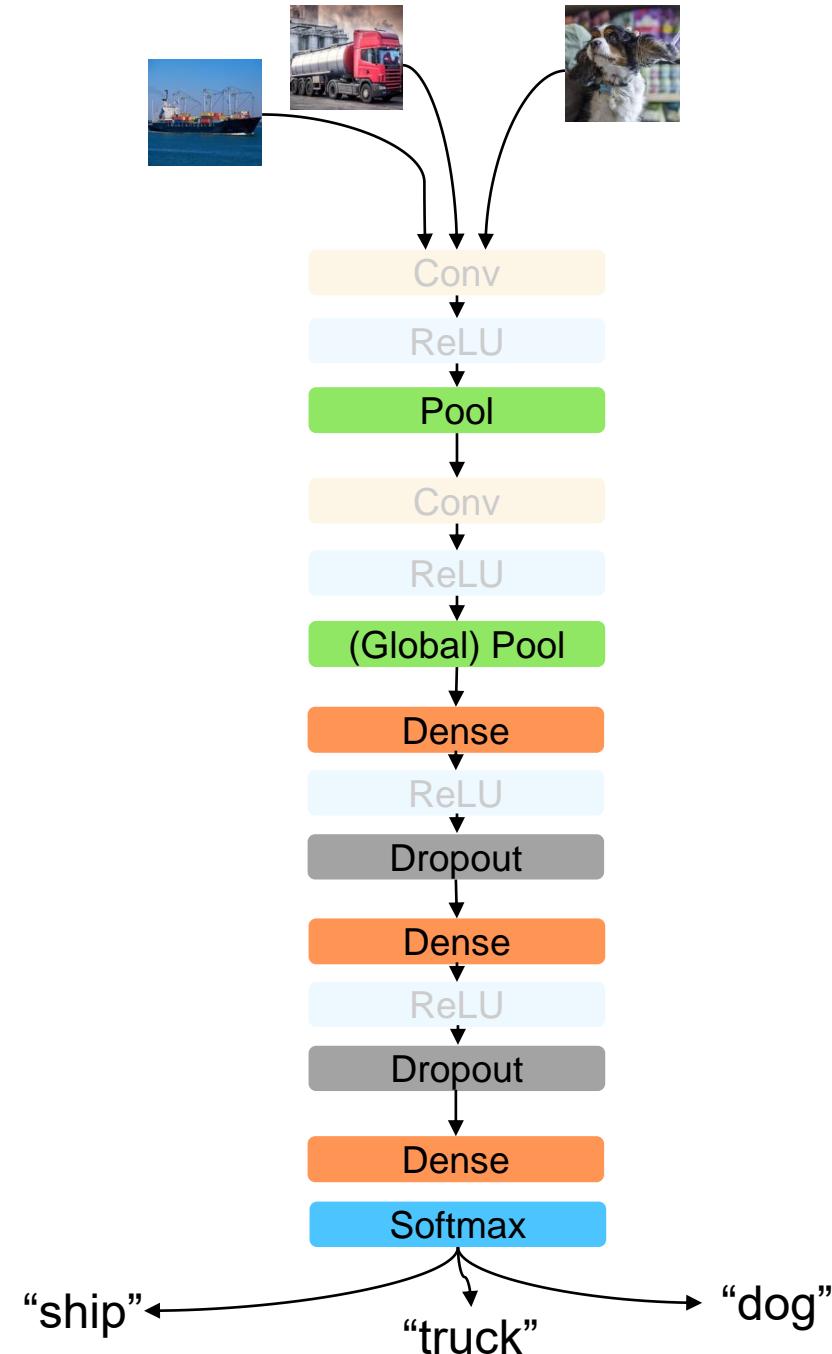


# CNN Architecture Part II

## Overview

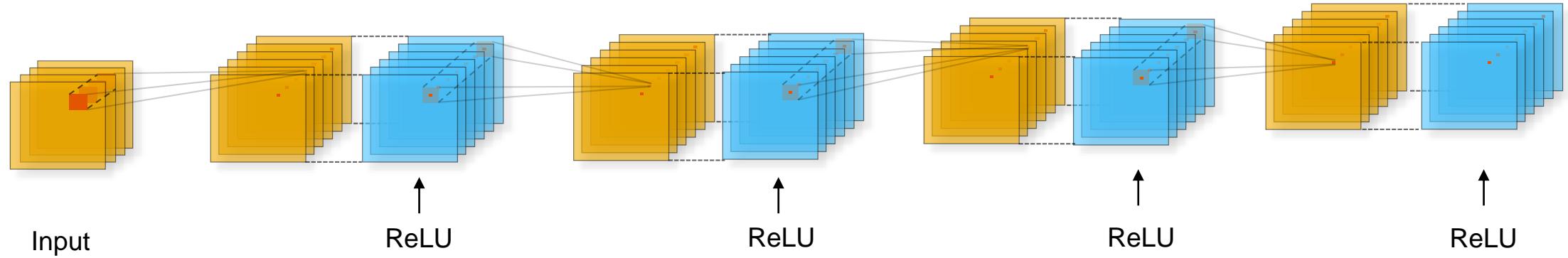
### Content:

- Convolutions
- ReLU non-linearity (activation function)
- Weight initialization
- Pooling / global pooling
- Dense (fully connected) layers
- Dropout
- Softmax



# CNN Architecture Part II

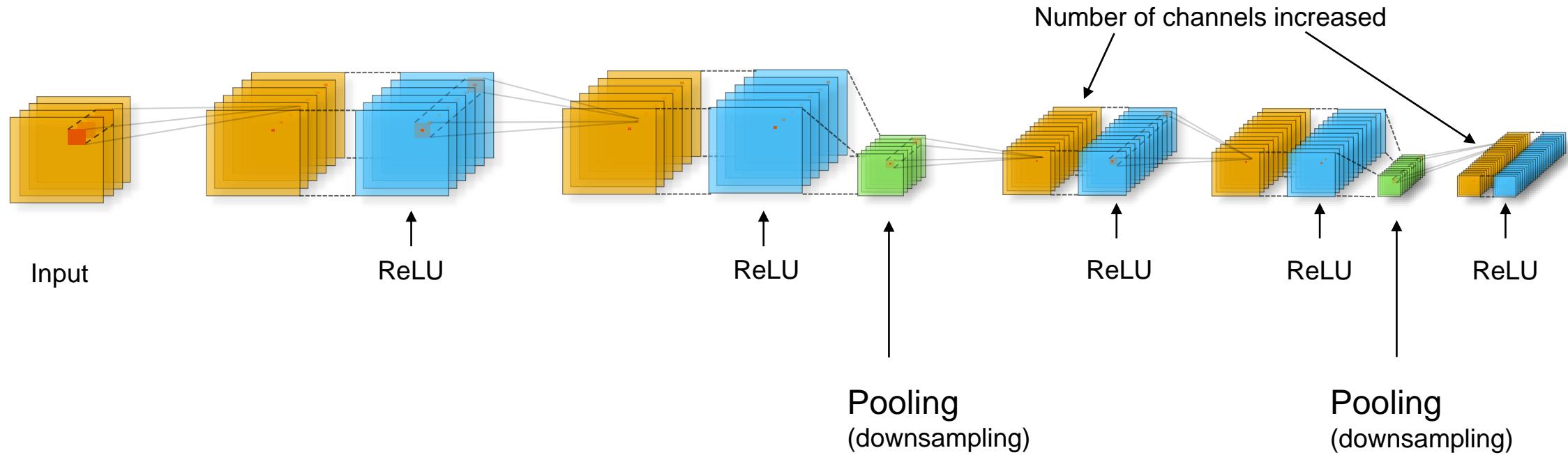
Reducing resolution: pooling



After each convolution (block), we may want to decrease resolution but increase the number of channels

# CNN Architecture Part II

Reducing resolution: pooling



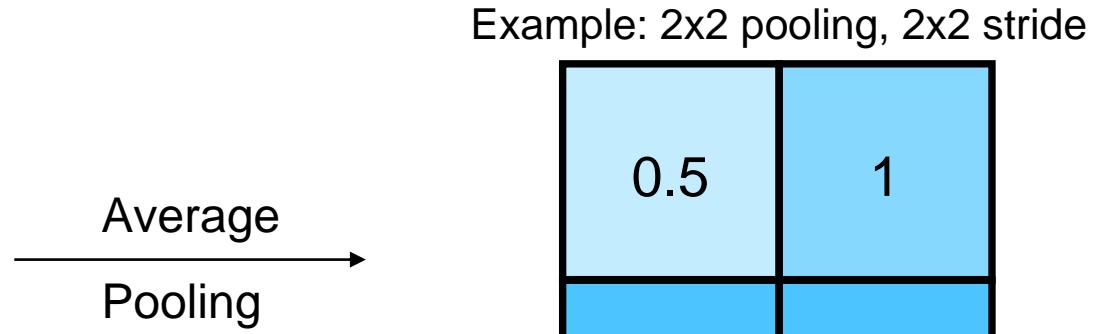
After each convolution (block), we may want to **decrease resolution** but **increase the number of channels**:

- Increase spatial invariance
- Increase level of abstraction
- Decrease computational load

# CNN Architecture Part II

Reducing resolution: pooling

0	1	2	1
1	0	0	1
2	3	1	0
1	2	5	2



Example: 2x2 pooling, 2x2 stride

0.5	1
2	2

# CNN Architecture Part II

Reducing resolution: pooling

0	1	2	1
1	0	0	1
2	3	1	0
1	2	5	2

Average  
Pooling →

Example: 2x2 pooling, 2x2 stride

0.5	1
2	2

0	1	2	1
1	0	0	1
2	3	1	0
1	2	5	2

Max  
Pooling →

1	2
3	5

Better at capturing  
edges

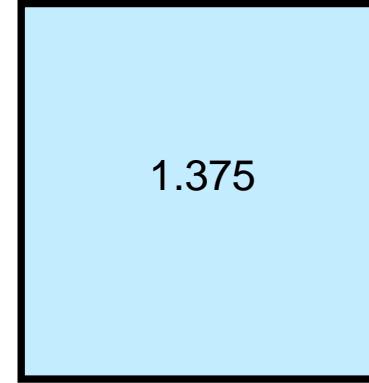
# CNN Architecture Part II

Reducing resolution: global pooling

0	1	2	1
1	0	0	1
2	3	1	0
1	2	5	2

Global Average  
Pooling

Take the average over an entire channel

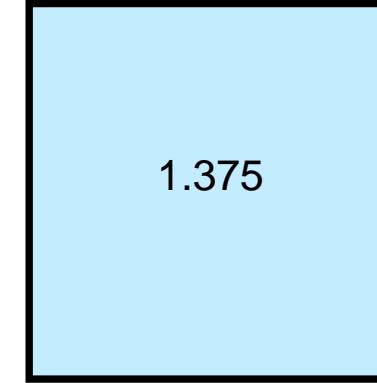


# CNN Architecture Part II

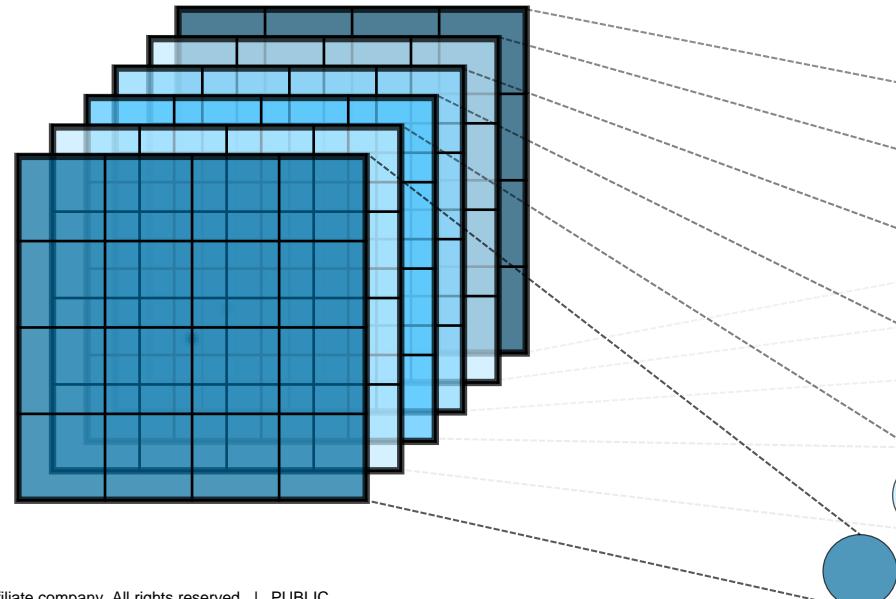
Reducing resolution: global pooling

0	1	2	1
1	0	0	1
2	3	1	0
1	2	5	2

Global Average  
Pooling



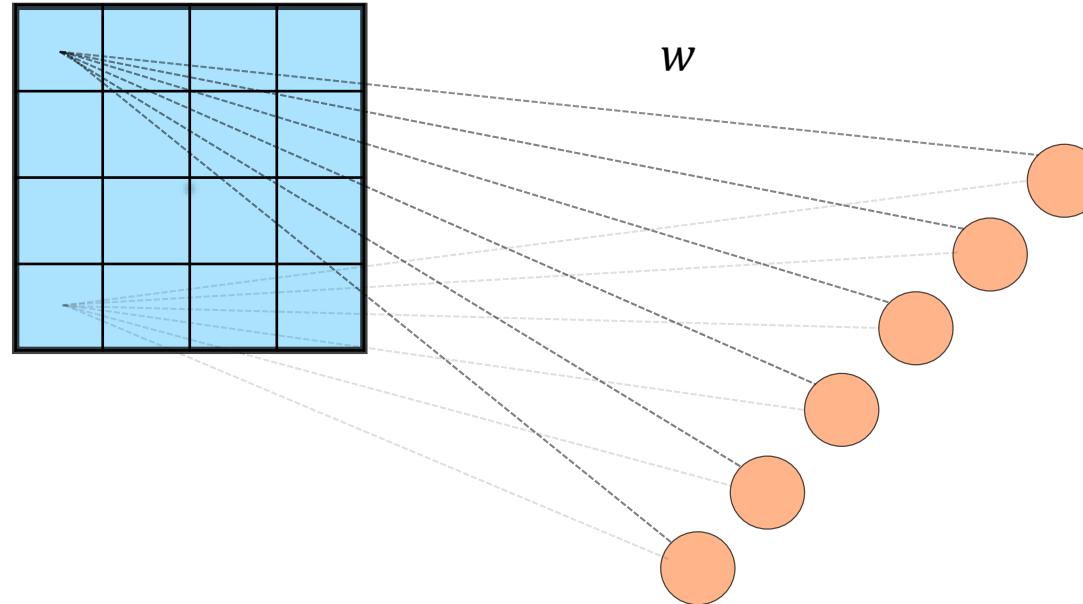
Take the average over an entire channel



Outcome of global pooling can be used as input for a dense layer to reduce complexity

# CNN Architecture Part II

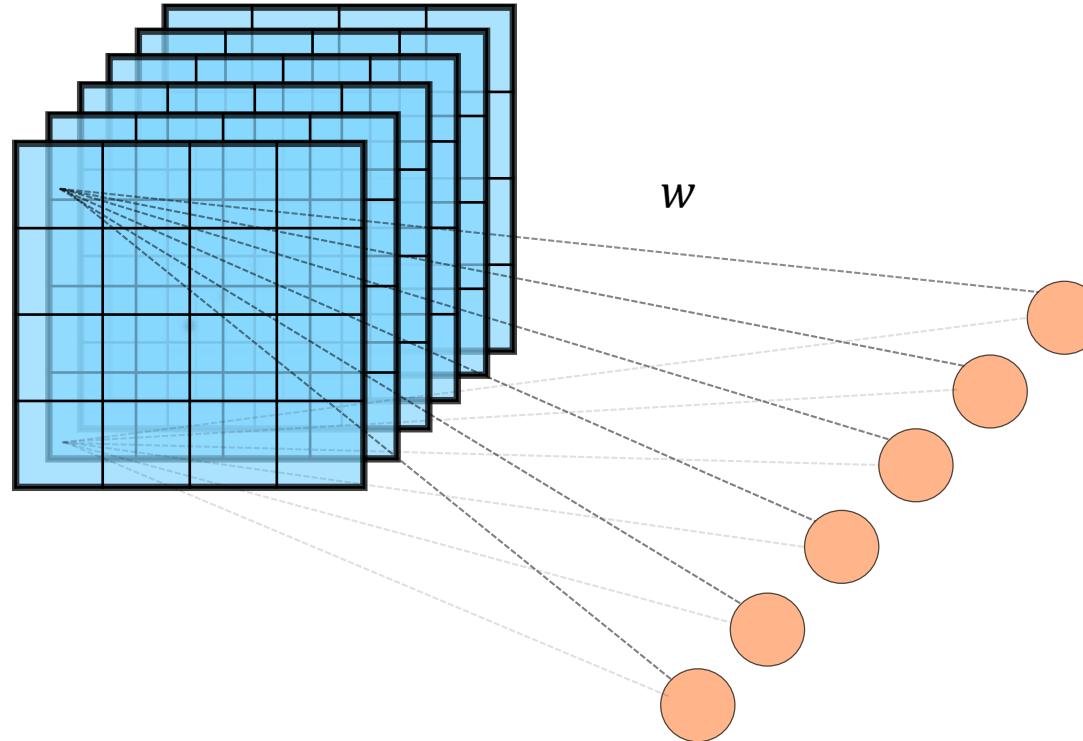
## From convolutions to dense layers



Every pixel in all output channels is connected to every neuron in the dense layer

# CNN Architecture Part II

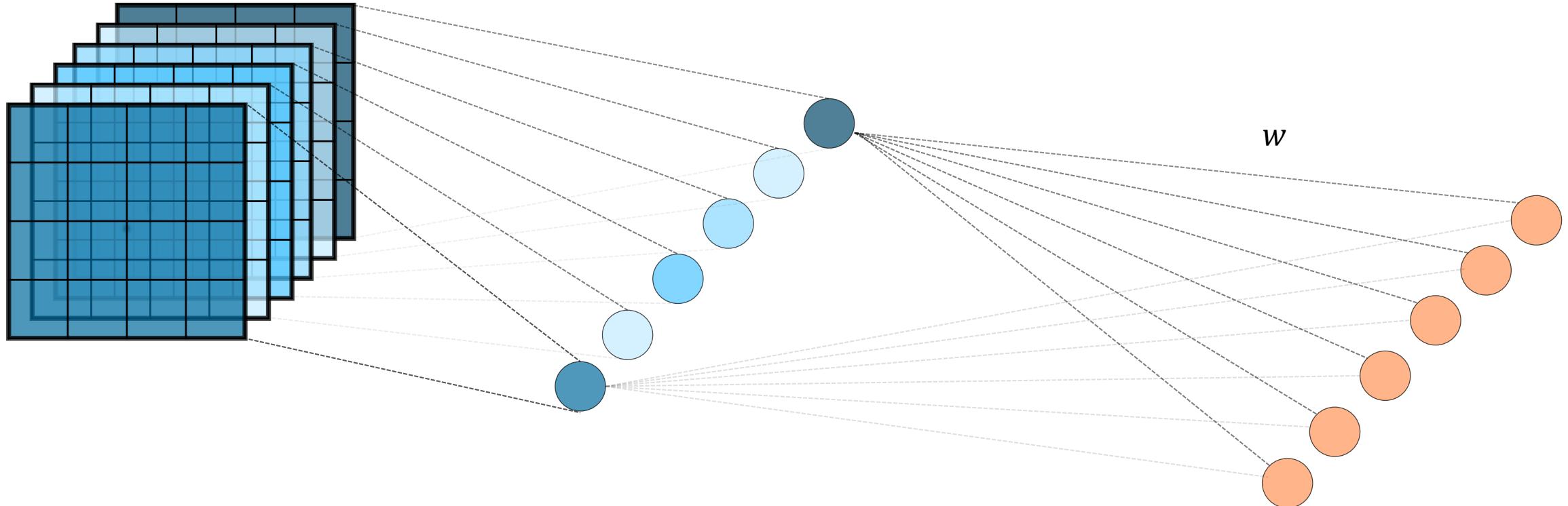
From convolutions to dense layers



Every pixel *in all output channels* is connected to every neuron in the dense layer  
→ weight matrix may become very large!

# CNN Architecture Part II

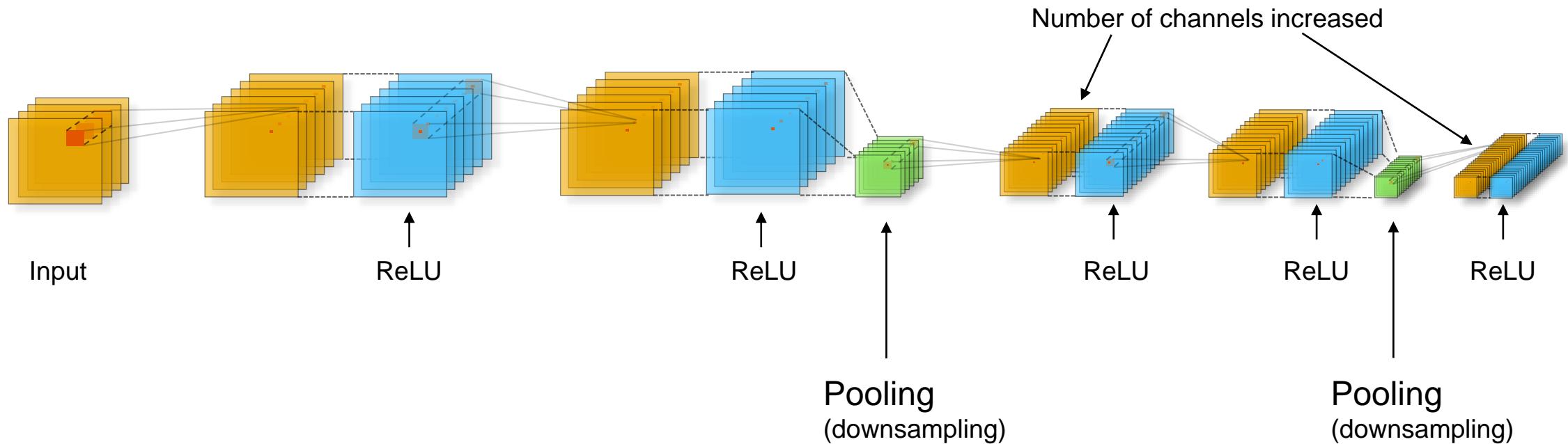
From convolutions to dense layers



Every pixel *in all output channels* is connected to every neuron in the dense layer  
→ weight matrix may become very large!  
→ use **global pooling**

# CNN Architecture Part II

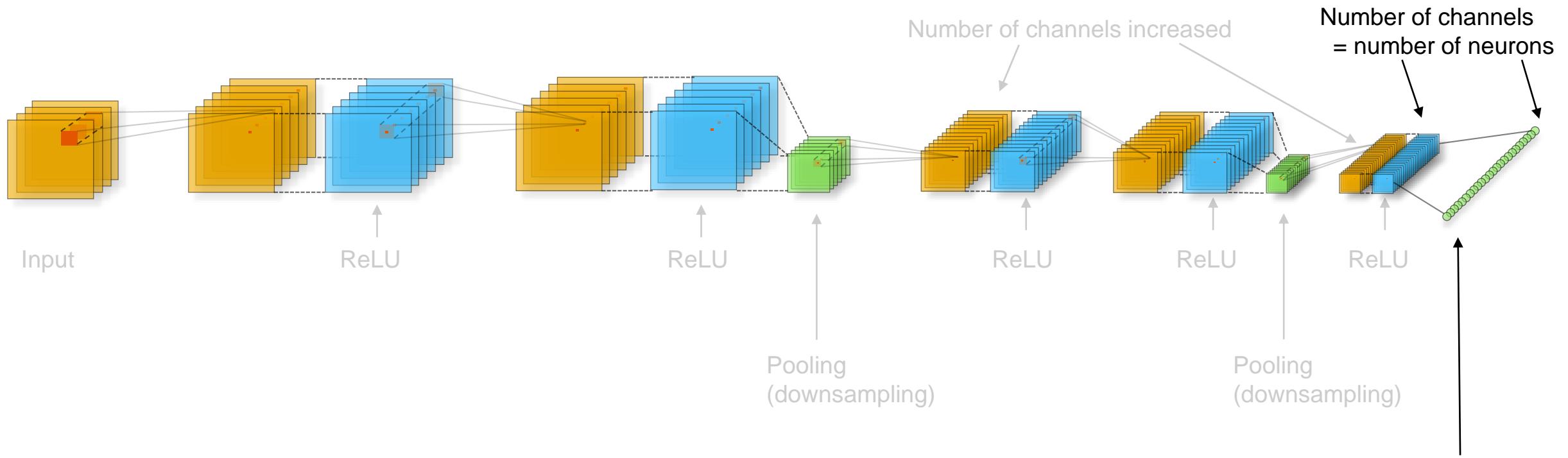
## From convolutions to dense layers



For classification tasks, we want to add dense layers

# CNN Architecture Part II

## From convolutions to dense layers

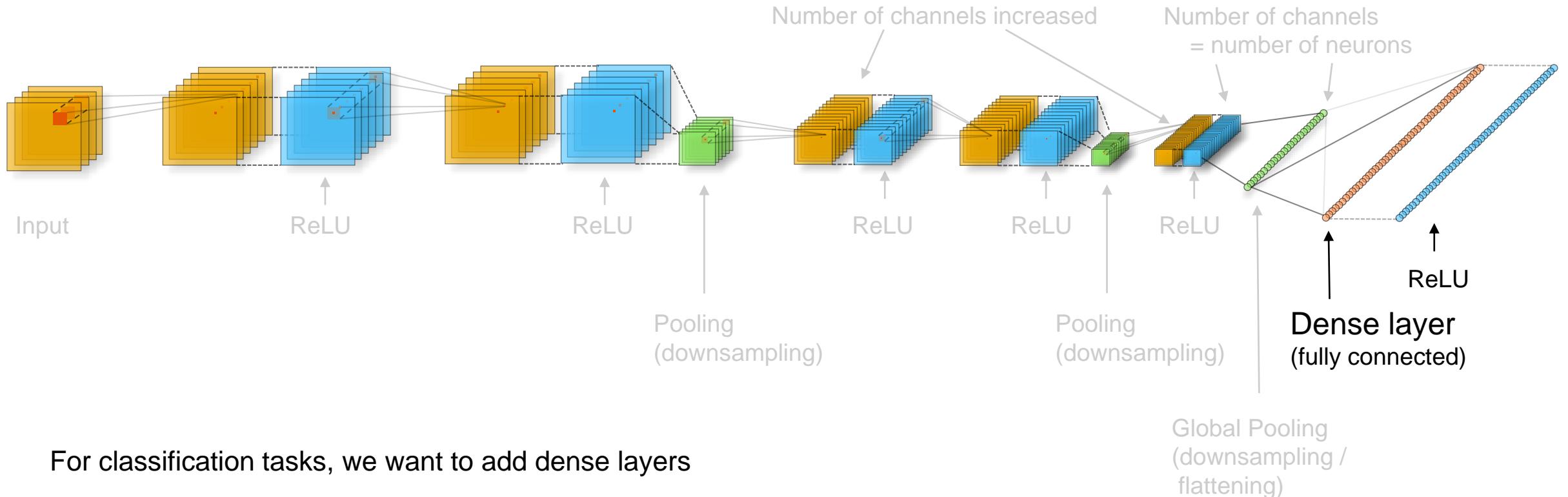


For classification tasks, we want to add dense layers

**Global Pooling**  
(downsampling /  
flattening)

# CNN Architecture Part II

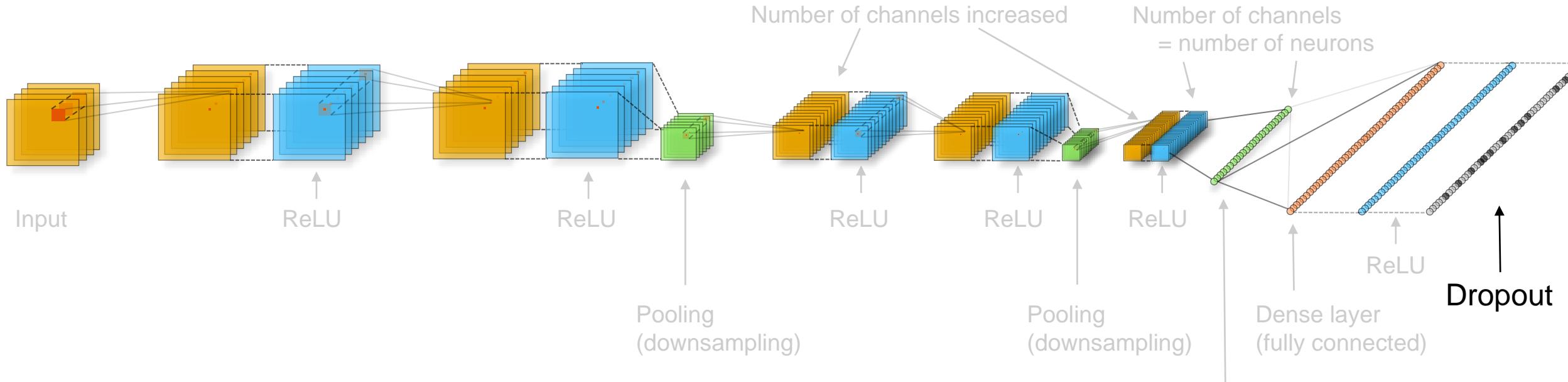
## From convolutions to dense layers



For classification tasks, we want to add dense layers

# CNN Architecture Part II

Risk of overfitting: dropout



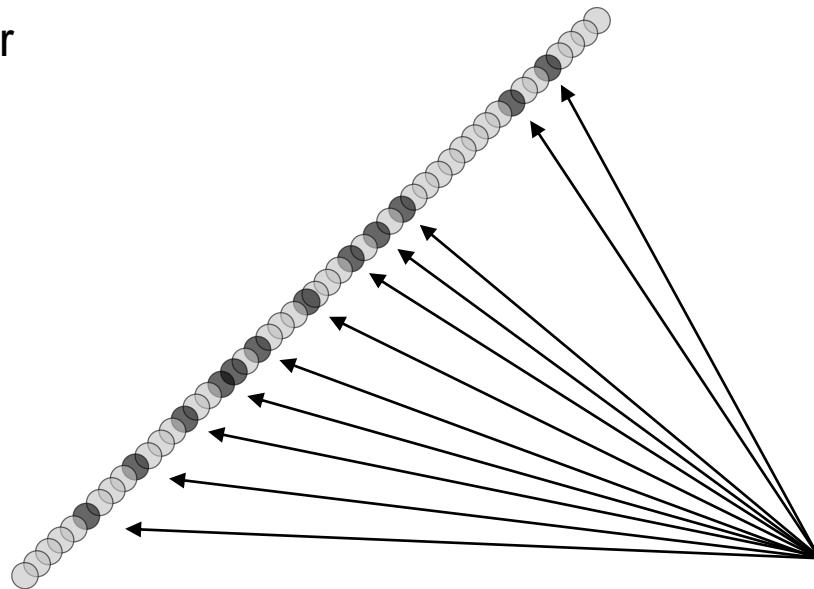
Dense layers have a lot of parameters and can therefore easily overfit

Remedy: Dropout

# CNN Architecture Part II

Risk of overfitting: dropout

Dropout Layer



## Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Srivastava, Hinton, Krizhevsky, Sutskever, Salakhutdinov  
Journal of Machine Learning Research 15 (2014) 1929-1958

*During training, randomly set neurons to zero  
with probability P (e.g. 50% chance)*

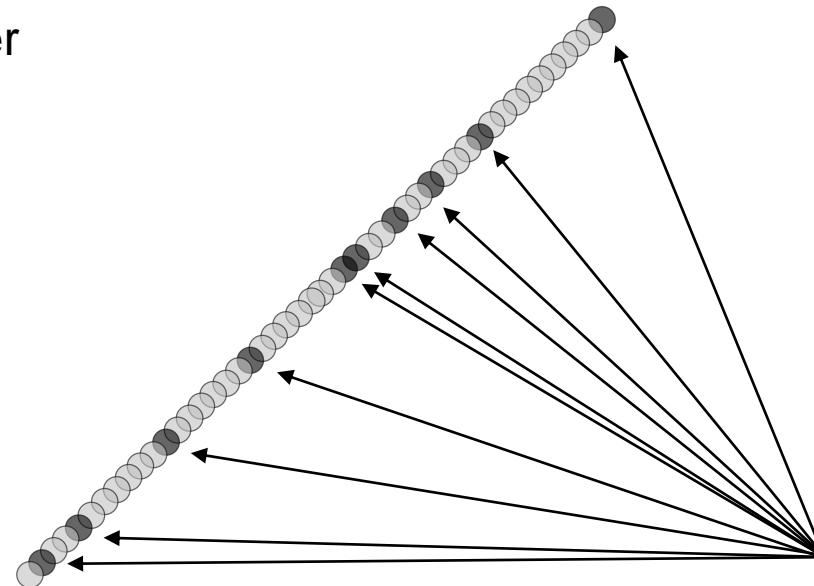
Dense layers have a lot of parameters and can therefore easily overfit

Remedy: Dropout

# CNN Architecture Part II

Risk of overfitting: dropout

Dropout Layer



## Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Srivastava, Hinton, Krizhevsky, Sutskever, Salakhutdinov  
Journal of Machine Learning Research 15 (2014) 1929-1958

*During training, randomly set neurons to zero  
with probability P (e.g. 50% chance)*

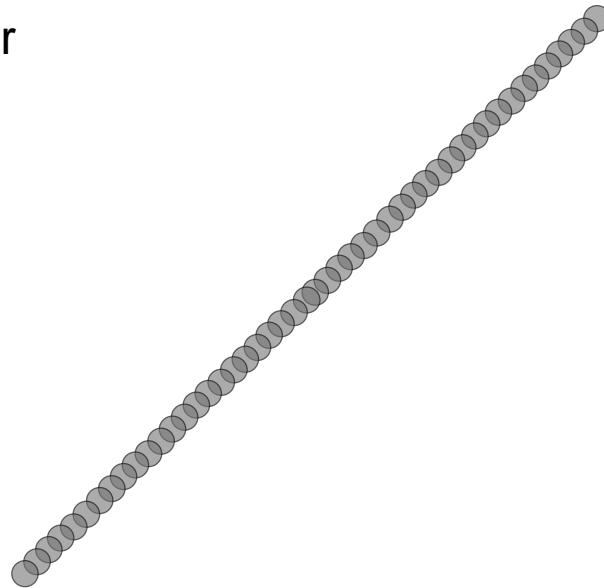
Dense layers have a lot of parameters and can therefore easily overfit

Remedy: Dropout

# CNN Architecture Part II

Risk of overfitting: dropout

Dropout Layer



**Dropout: A Simple Way to Prevent Neural Networks from Overfitting**

Srivastava, Hinton, Krizhevsky, Sutskever, Salakhutdinov  
Journal of Machine Learning Research 15 (2014) 1929-1958

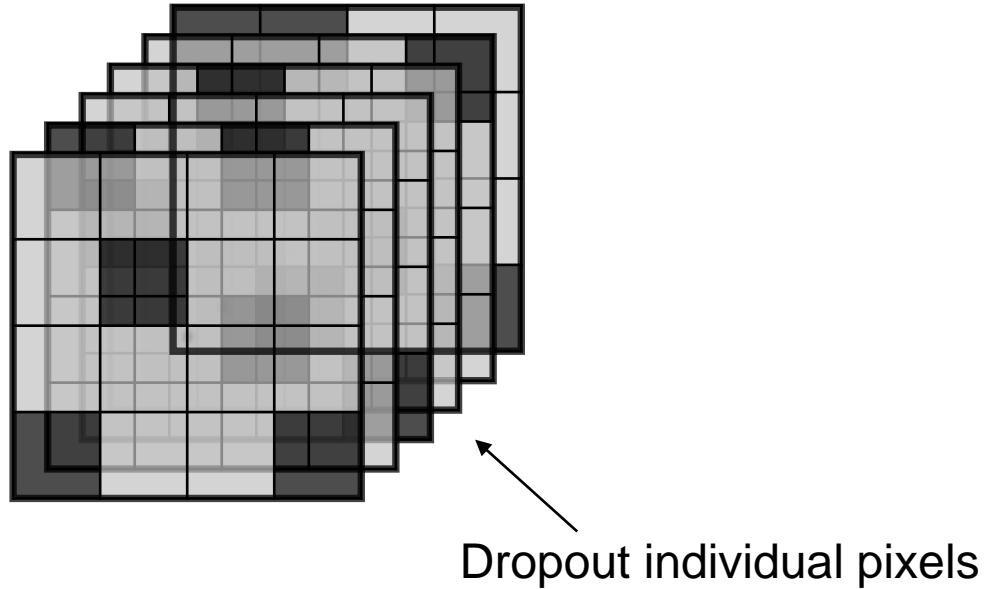
*During inference, rescale neurons  
with  $1 / P$  (e.g. factor 2)*

Dense layers have a lot of parameters and can therefore easily overfit

Remedy: **Dropout**

# CNN Architecture Part II

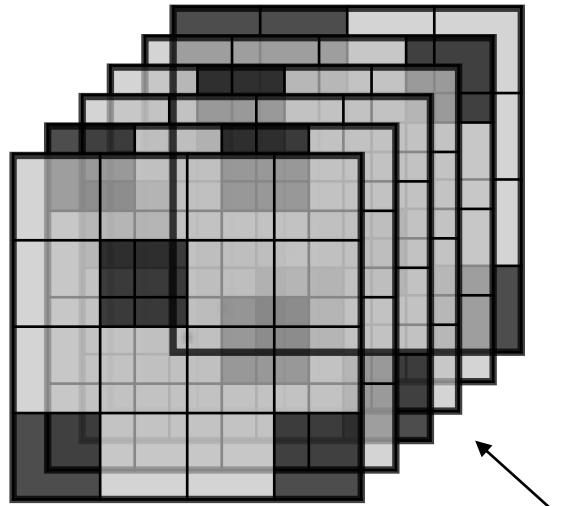
Risk of overfitting: dropout



**Dropout** can also be applied to **convolutional layers**

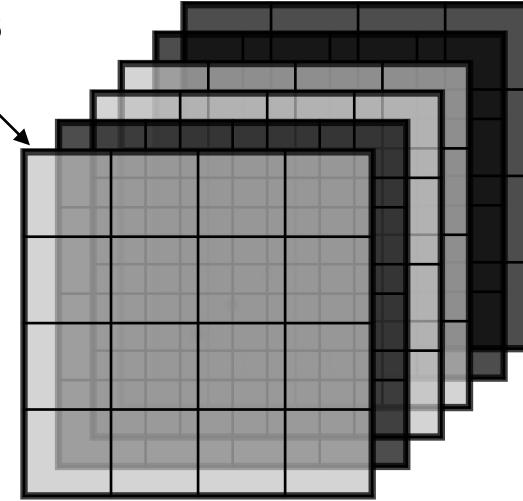
# CNN Architecture Part II

Risk of overfitting: dropout



Dropout individual pixels

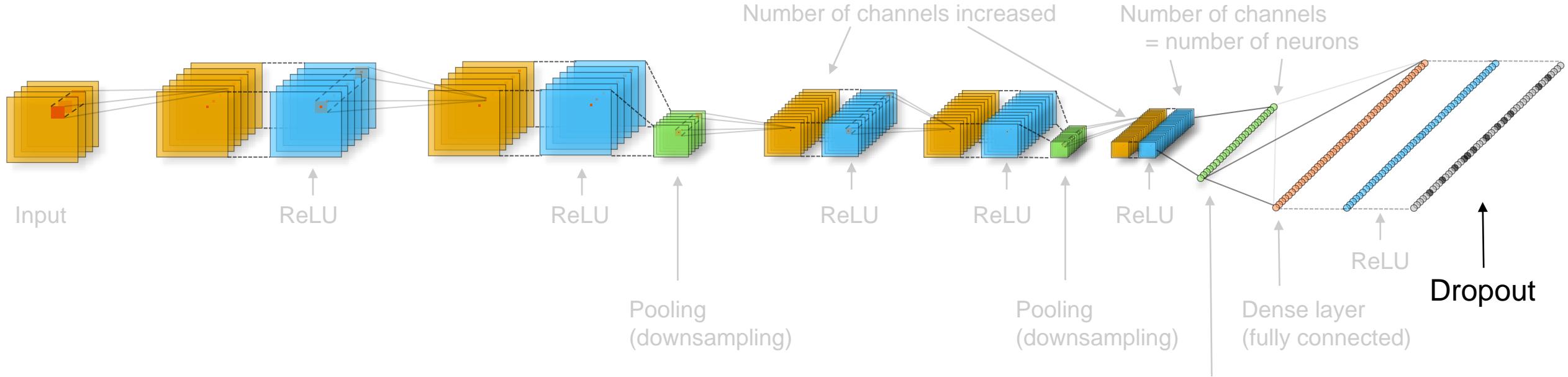
Dropout entire channels  
(aka spatial dropout)



Dropout can also be applied to **convolutional layers**

# CNN Architecture Part II

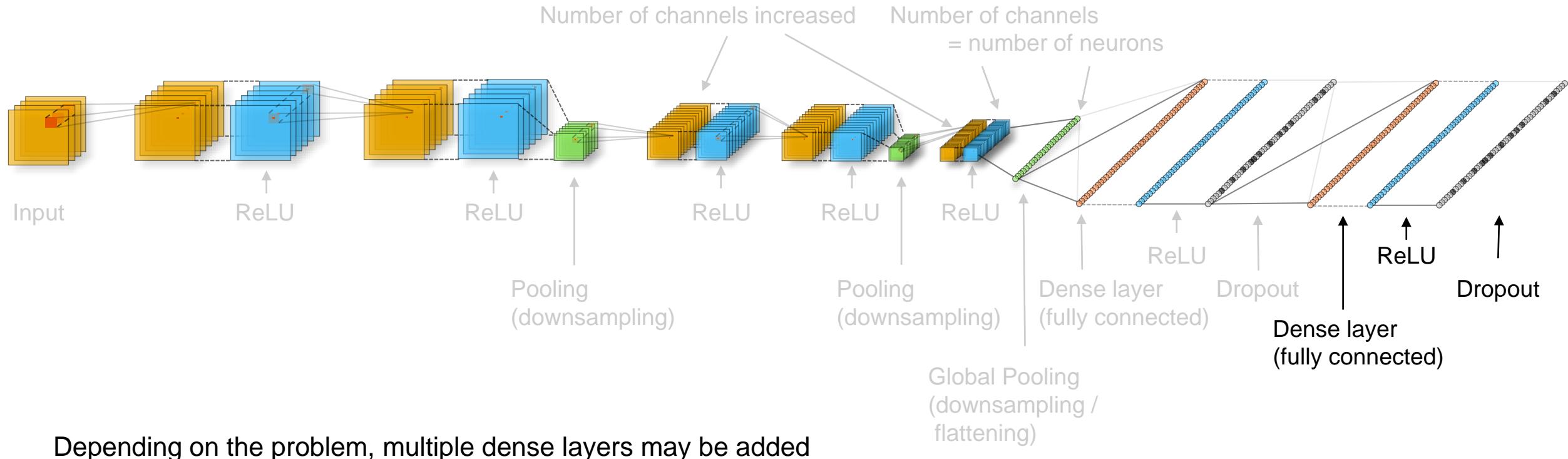
## Multiple dense layers



Depending on the problem, multiple dense layers may be added

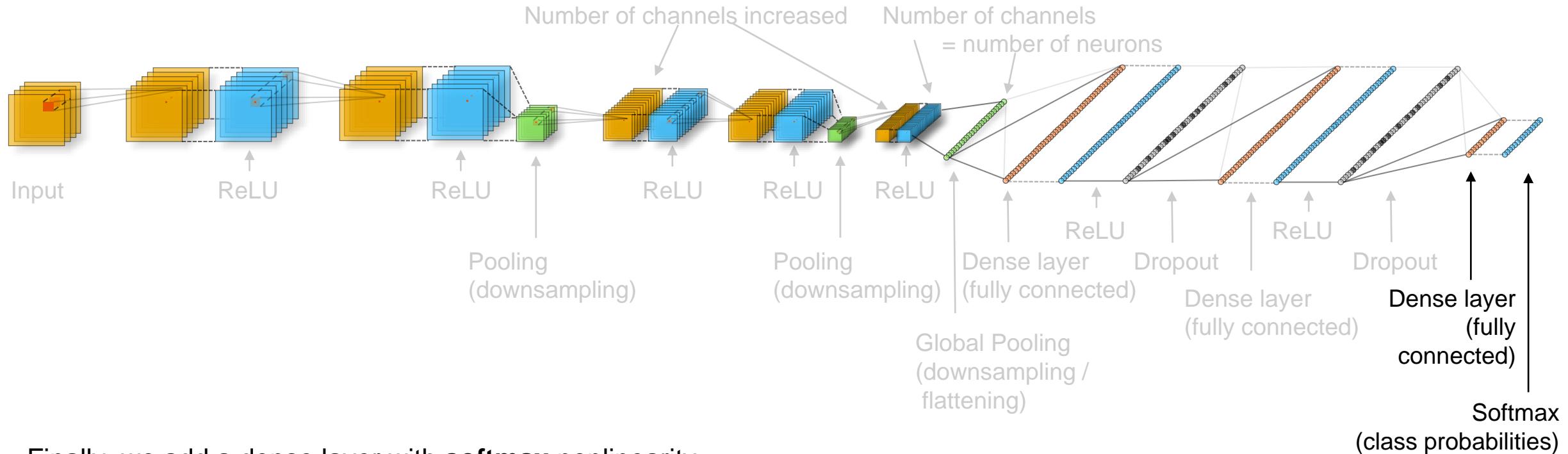
# CNN Architecture Part II

## Multiple dense layers



# CNN Architecture Part II

## Final dense & classification layer

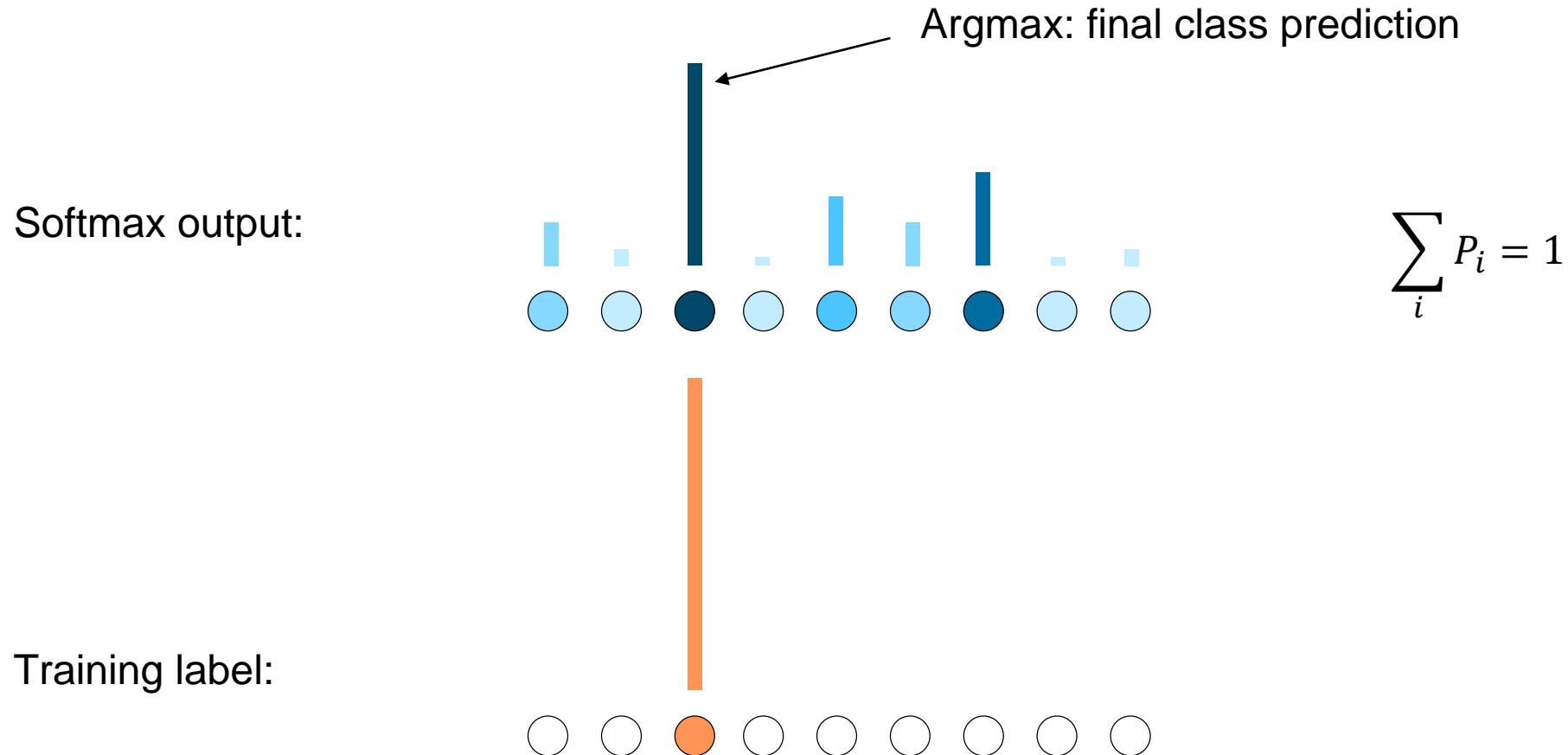


Finally, we add a dense layer with **softmax** nonlinearity

The number of neurons in the final dense layer must equal the number of classes

# CNN Architecture Part II

Final dense & classification layer

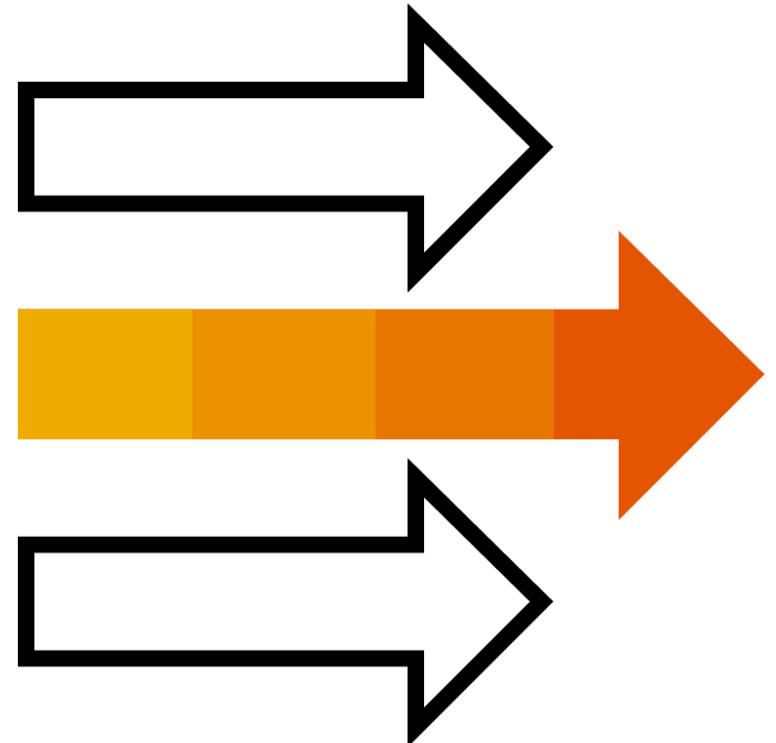


# CNN Architecture Part II

Coming up next

## Accelerating Deep CNN Training

- Computational considerations
- Batch normalization
- Transfer learning
- Residual networks



# Thank you.

Contact information:

[open@sap.com](mailto:open@sap.com)

# © 2017 SAP SE or an SAP affiliate company. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP SE or an SAP affiliate company.

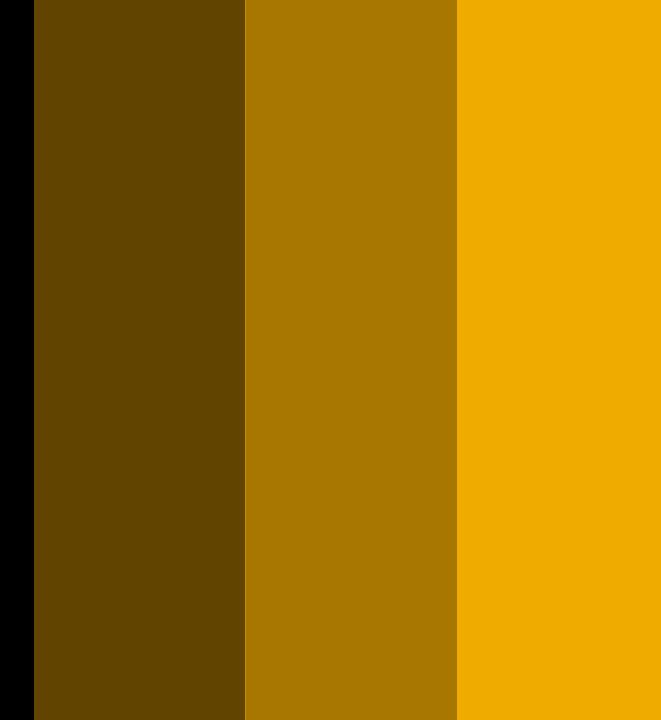
The information contained herein may be changed without prior notice. Some software products marketed by SAP SE and its distributors contain proprietary software components of other software vendors. National product specifications may vary.

These materials are provided by SAP SE or an SAP affiliate company for informational purposes only, without representation or warranty of any kind, and SAP or its affiliated companies shall not be liable for errors or omissions with respect to the materials. The only warranties for SAP or SAP affiliate company products and services are those that are set forth in the express warranty statements accompanying such products and services, if any. Nothing herein should be construed as constituting an additional warranty.

In particular, SAP SE or its affiliated companies have no obligation to pursue any course of business outlined in this document or any related presentation, or to develop or release any functionality mentioned therein. This document, or any related presentation, and SAP SE's or its affiliated companies' strategy and possible future developments, products, and/or platform directions and functionality are all subject to change and may be changed by SAP SE or its affiliated companies at any time for any reason without notice. The information in this document is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. All forward-looking statements are subject to various risks and uncertainties that could cause actual results to differ materially from expectations. Readers are cautioned not to place undue reliance on these forward-looking statements, and they should not be relied upon in making purchasing decisions.

SAP and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP SE (or an SAP affiliate company) in Germany and other countries. All other product and service names mentioned are the trademarks of their respective companies.

See <http://global.sap.com/corporate-en/legal/copyright/index.epx> for additional trademark information and notices.



Week 4: Convolutional Networks

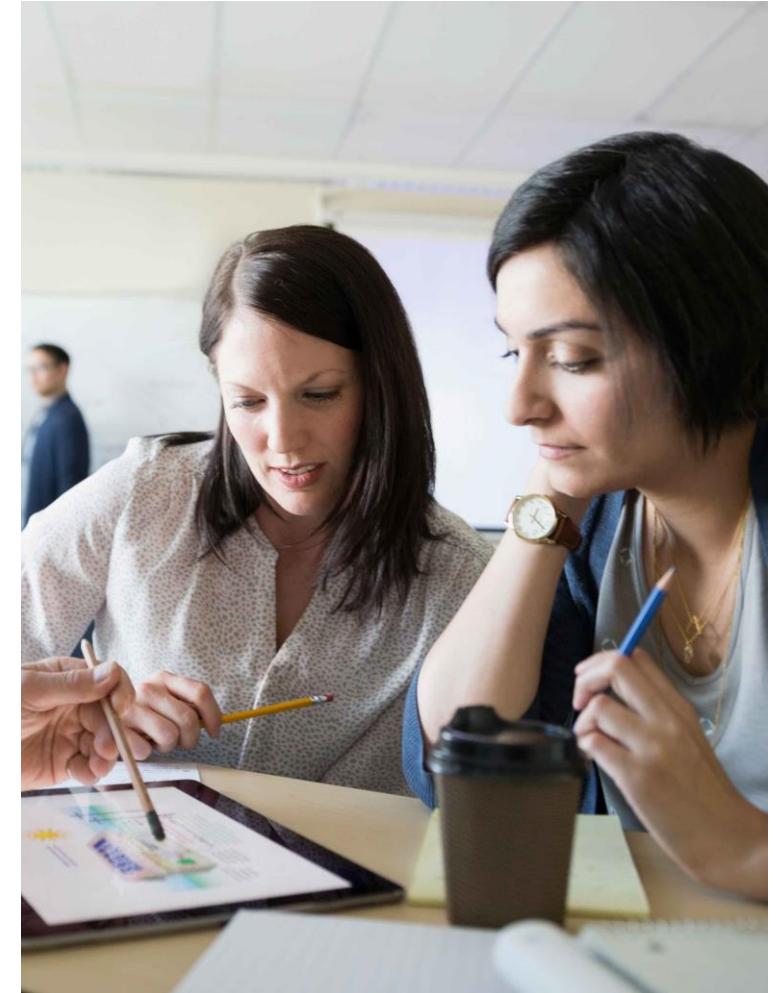
## **Unit 4: Accelerating Deep CNN Training**

# Accelerating Deep CNN Training

What we covered in the last unit

## CNN Architecture Part II

- Pooling
- Dense layers
- Dropout
- Softmax

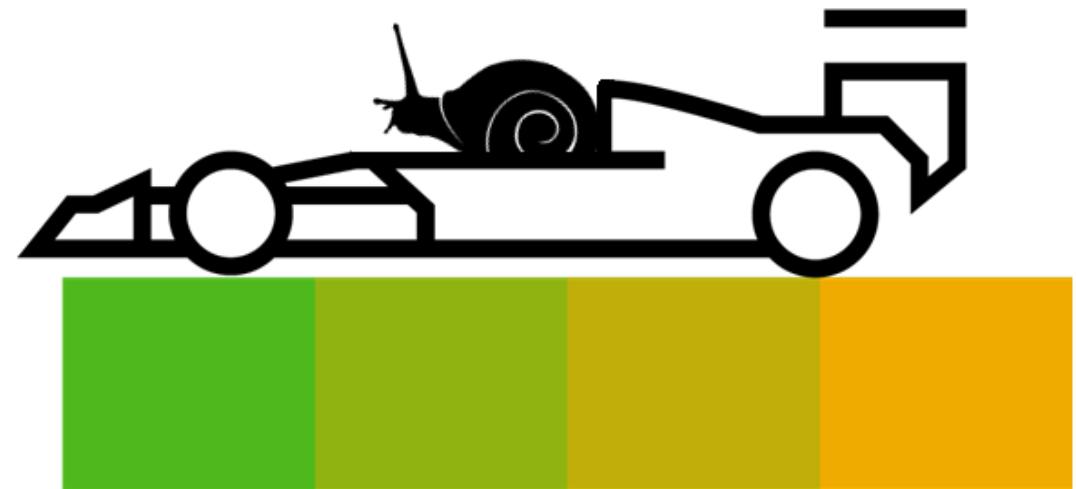


# Accelerating Deep CNN Training

## Overview

### Content:

- Computational considerations
- Batch normalization
- Transfer learning
- Residual networks



# Accelerating Deep CNN Training

## Overview

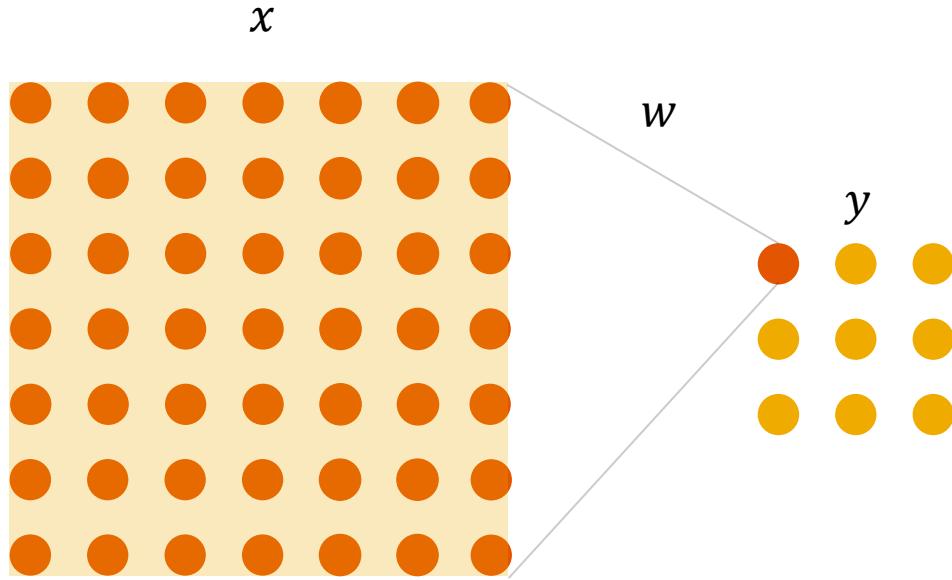
### Content:

- Computational considerations
- Batch normalization
- Transfer learning
- Residual networks

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

# Accelerating Deep CNN Training

Computational considerations: convolutional filter size



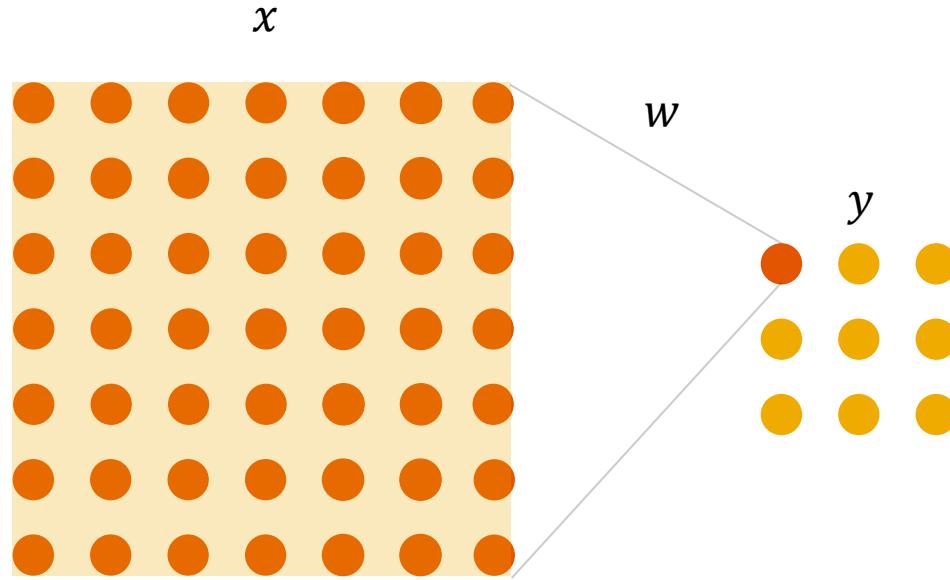
## 7 x 7 Convolution

100 channels on each layer

$49 \times 100^2$  parameters to optimize

# Accelerating Deep CNN Training

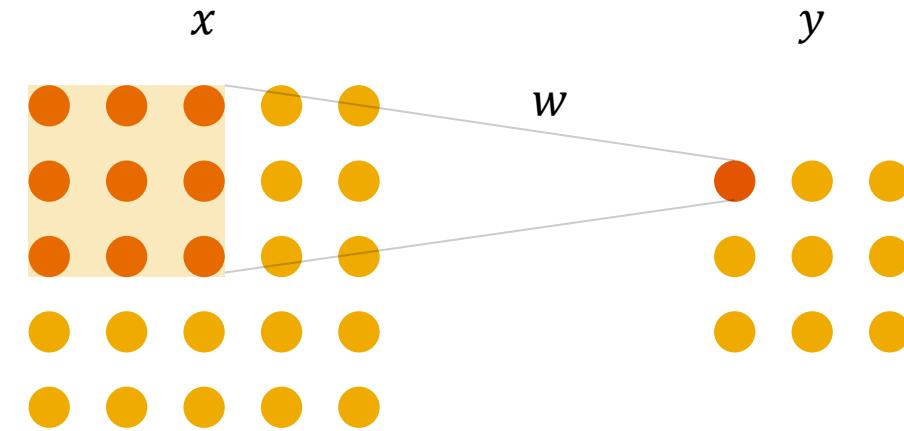
Computational considerations: convolutional filter size



## 7 x 7 Convolution

100 channels on each layer

$49 \times 100^2$  parameters to optimize



## 3 x 3 Convolution

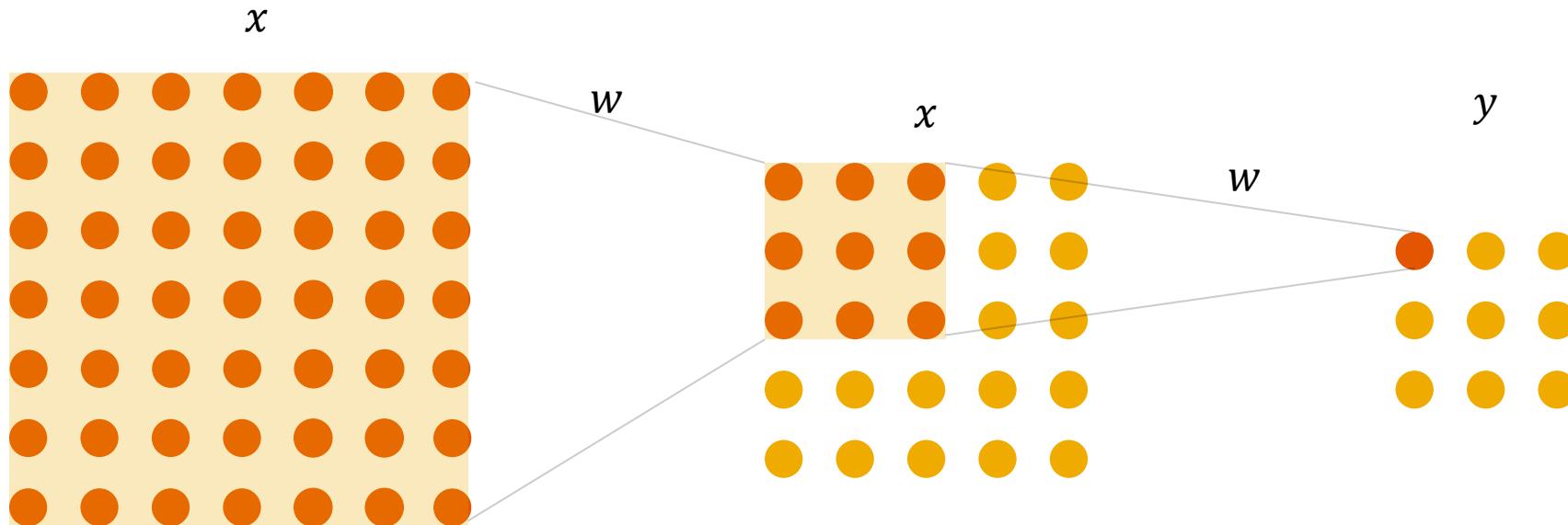
100 channels on each layer

$9 \times 100^2$  parameters to optimize

80% less

# Accelerating Deep CNN Training

Computational considerations: convolutional filter size



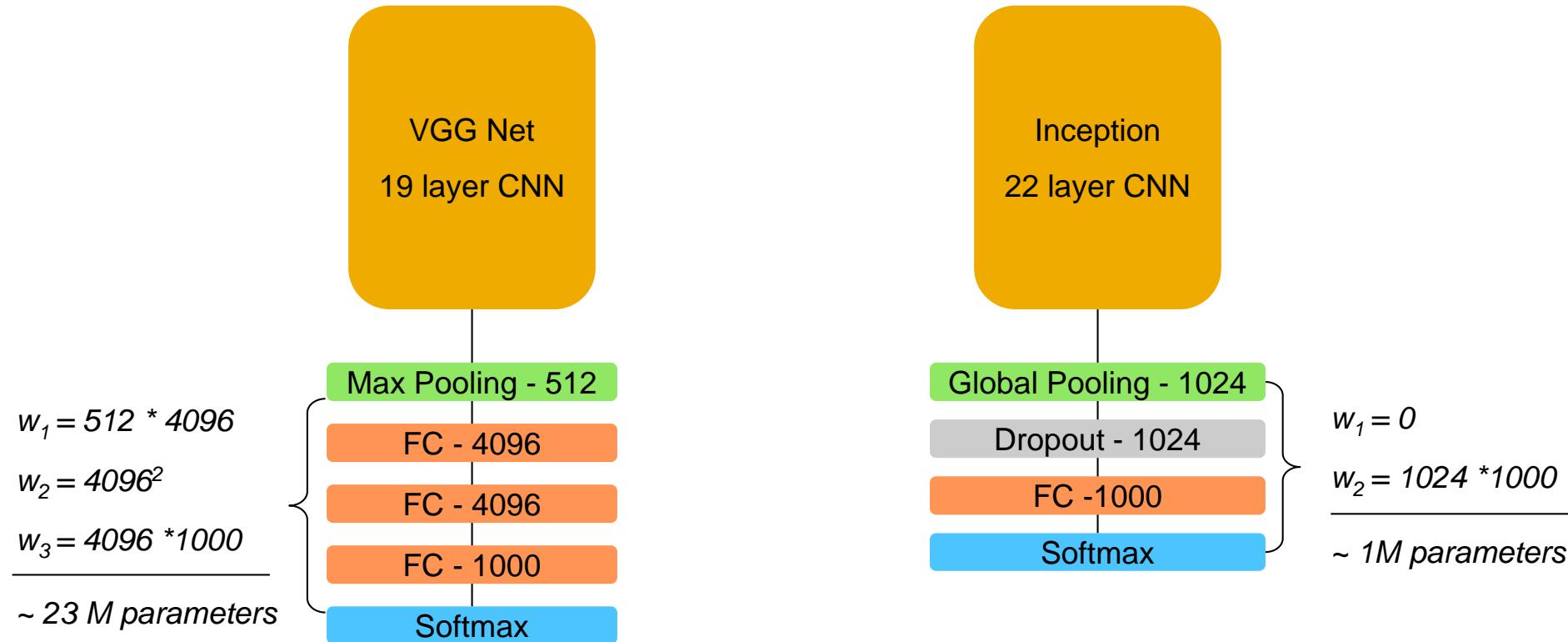
Combination of two 3x3 convolutional layers (2 strides) have the same field of view.

Still **fewer** parameters than 7 x 7, **less** overfitting, **more** nonlinearity!

**VGG Net:** Simonyan, K. & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR, abs/1409.1556.

# Accelerating Deep CNN Training

Computational considerations: replacing fully connected (dense) layers



**VGG Net:** Simonyan, K. & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR, abs/1409.1556.

**Inception:** Szegedy C. et al. (2015). Going Deeper with Convolutions. CoRR, abs/1409.4842.

# Accelerating Deep CNN Training

## Overview

### Content:

- Computational considerations
- **Batch normalization**
- Transfer learning
- Residual networks

$$\sqrt{b^2 - 4ac}$$

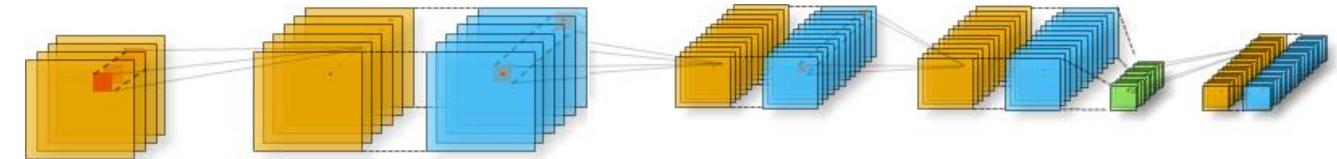
$$\int x dy$$

$$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$

$$e^{-i\omega t}$$

# Accelerating Deep CNN Training

Batch normalization: motivation

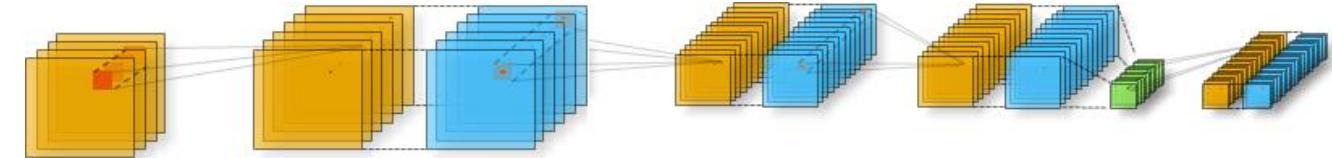


- During training, the output distribution of each layer changes due to weight updates

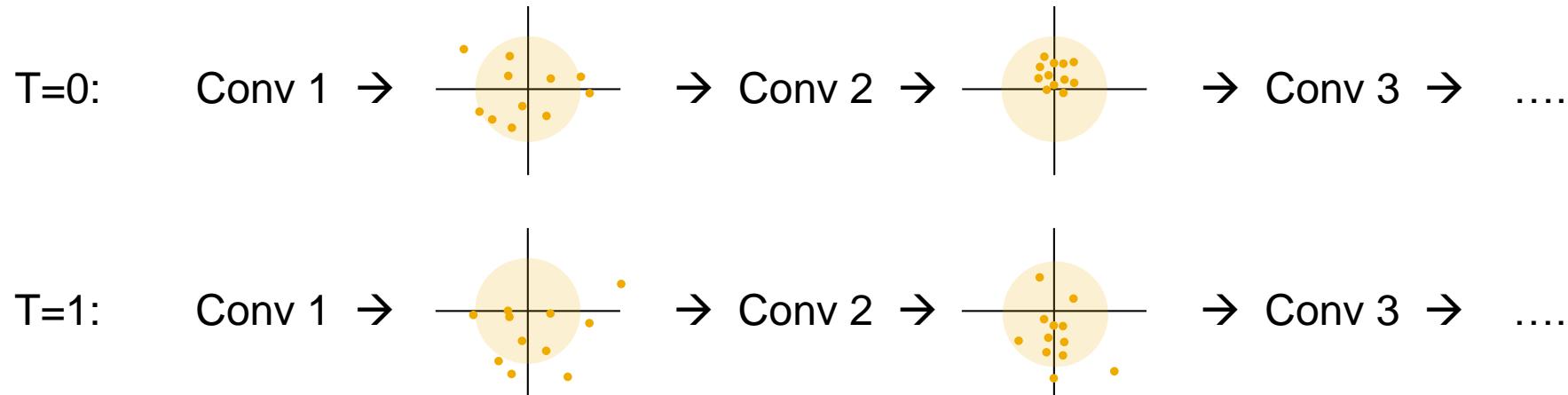


# Accelerating Deep CNN Training

Batch normalization: motivation

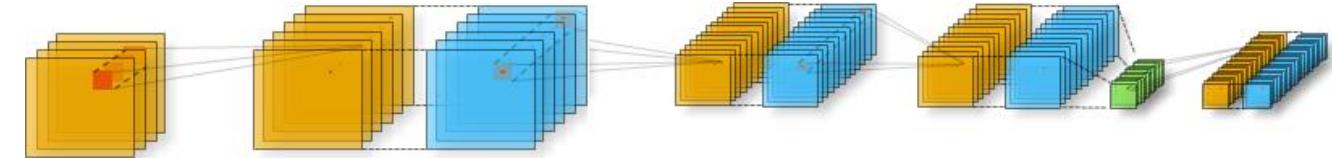


- During training, the output distribution of each layer changes due to weight updates

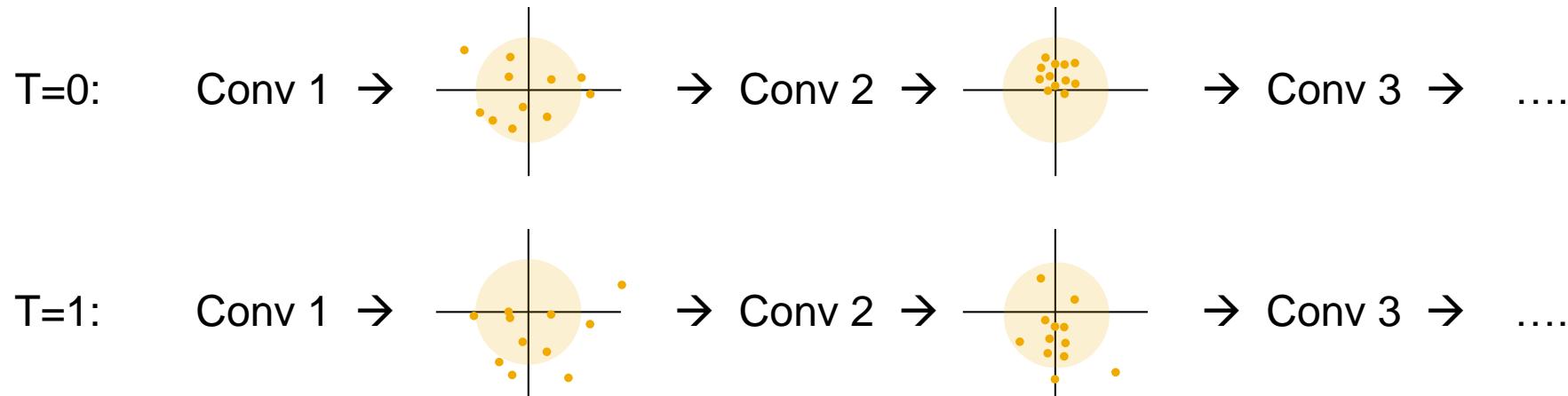


# Accelerating Deep CNN Training

Batch normalization: motivation



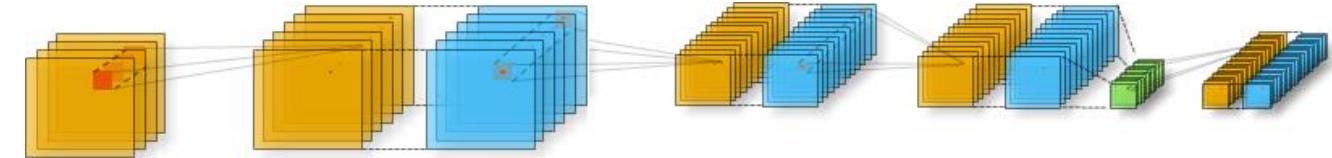
- During training, the output distribution of each layer changes due to weight updates
- Each following layer has to adapt to the new output distribution of the previous layer
- This makes training rather hard



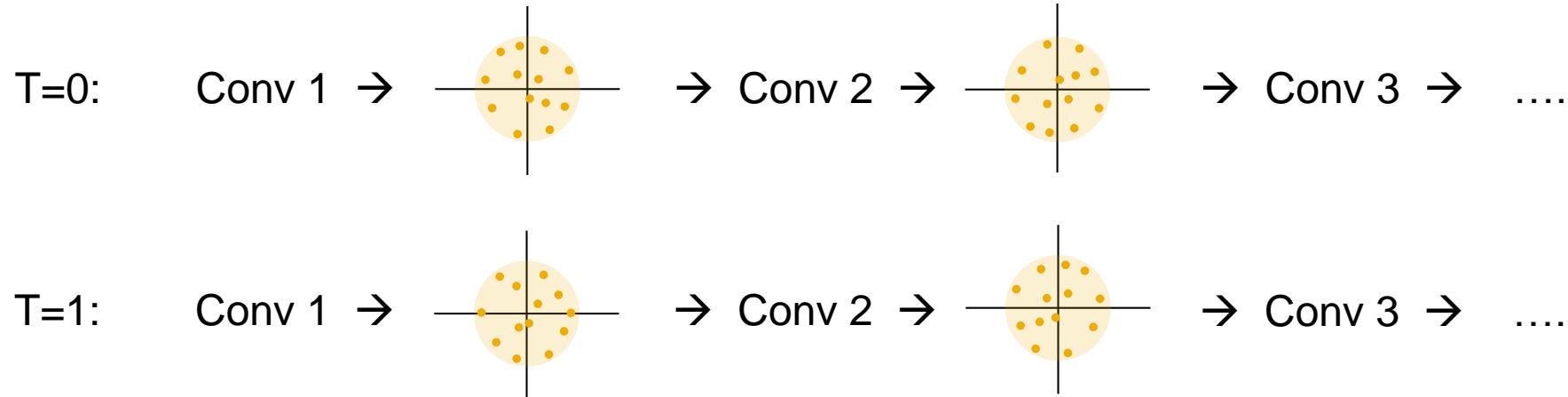
Aka “internal covariate shift”

# Accelerating Deep CNN Training

Batch normalization: motivation



- Idea: Normalize output distribution after each layer



# Accelerating Deep CNN Training

Batch normalization algorithm

**Input:** Values of  $x$  over a mini-batch  $B = \{x_{i..m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

**Batch Normalization:** Ioffe, S. & Szegedy C. (2015). *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. CoRR, abs/1502.03167.

# Accelerating Deep CNN Training

Batch normalization algorithm

**Input:** Values of  $x$  over a mini-batch  $B = \{x_{i..m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

1. **Normalize every batch by mean and variance of the batch.**

$$\hat{x}_i \leftarrow \frac{x_i - \mu_\beta}{\sqrt{\alpha_\beta^2 + \epsilon}}$$

**Batch Normalization:** Ioffe, S. & Szegedy C. (2015). *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. CoRR, abs/1502.03167.

# Accelerating Deep CNN Training

Batch normalization algorithm

**Input:** Values of  $x$  over a mini-batch  $B = \{x_{i..m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

1. **Normalize every batch by mean and variance** of the batch.

$$\hat{x}_i \leftarrow \frac{x_i - \mu_\beta}{\sqrt{\alpha_\beta^2 + \epsilon}}$$

2. Introduce **two new parameters**  $\gamma, \beta$ .

**Batch Normalization:** Ioffe, S. & Szegedy C. (2015). *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. CoRR, abs/1502.03167.

# Accelerating Deep CNN Training

Batch normalization algorithm

**Input:** Values of  $x$  over a mini-batch  $B = \{x_{i..m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

1. **Normalize every batch by mean and variance** of the batch.

$$\hat{x}_i \leftarrow \frac{x_i - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}}$$

2. Introduce **two new parameters**  $\gamma, \beta$ .

3. **Scale and shift** the activations with  $\gamma$  and  $\beta$ .

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i)$$

**Batch Normalization:** Ioffe, S. & Szegedy C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. CoRR, abs/1502.03167.

# Accelerating Deep CNN Training

Batch normalization algorithm

**Input:** Values of  $x$  over a mini-batch  $B = \{x_{i..m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

1. **Normalize every batch by mean and variance** of the batch.

$$\hat{x}_i \leftarrow \frac{x_i - \mu_\beta}{\sqrt{\alpha_\beta^2 + \epsilon}}$$

2. Introduce **two new parameters**  $\gamma, \beta$ .

3. **Scale and shift** the activations with  $\gamma$  and  $\beta$ .

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i)$$

**Batch Normalization:** Ioffe, S. & Szegedy C. (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. CoRR, abs/1502.03167.

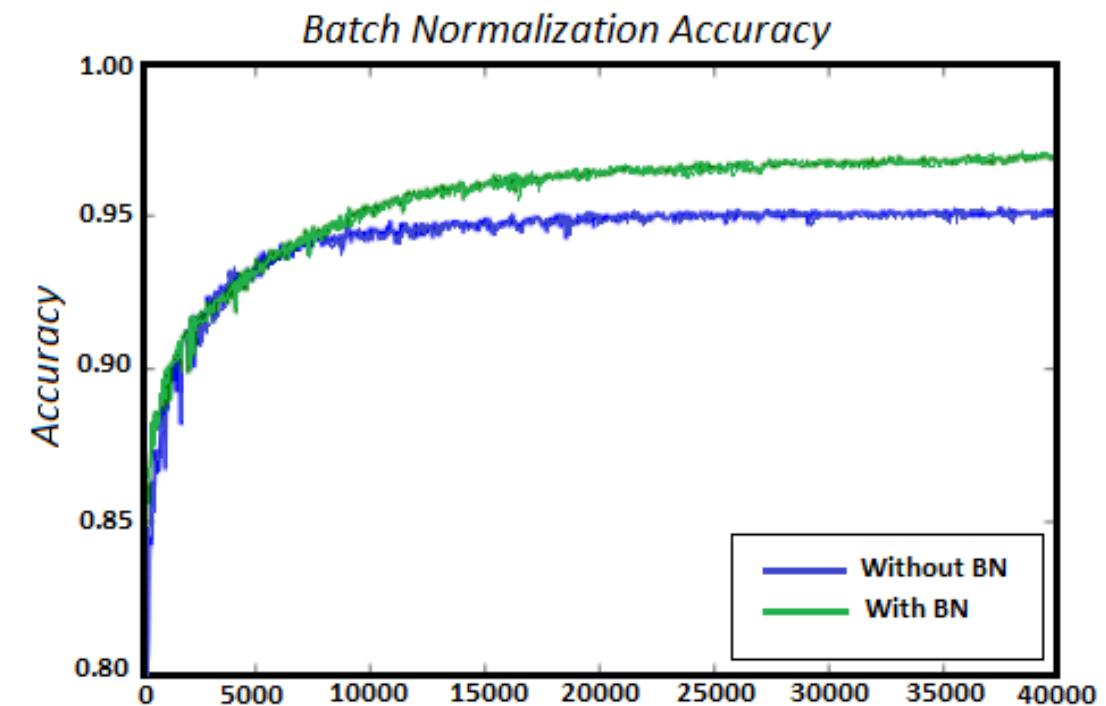
4. **Learn** the parameters  $\gamma$  and  $\beta$  during training.

# Accelerating Deep CNN Training

Batch normalization algorithm

## Benefits:

- Higher accuracy in earlier training steps
- Reduction in overfitting
- Less need for dropout layers

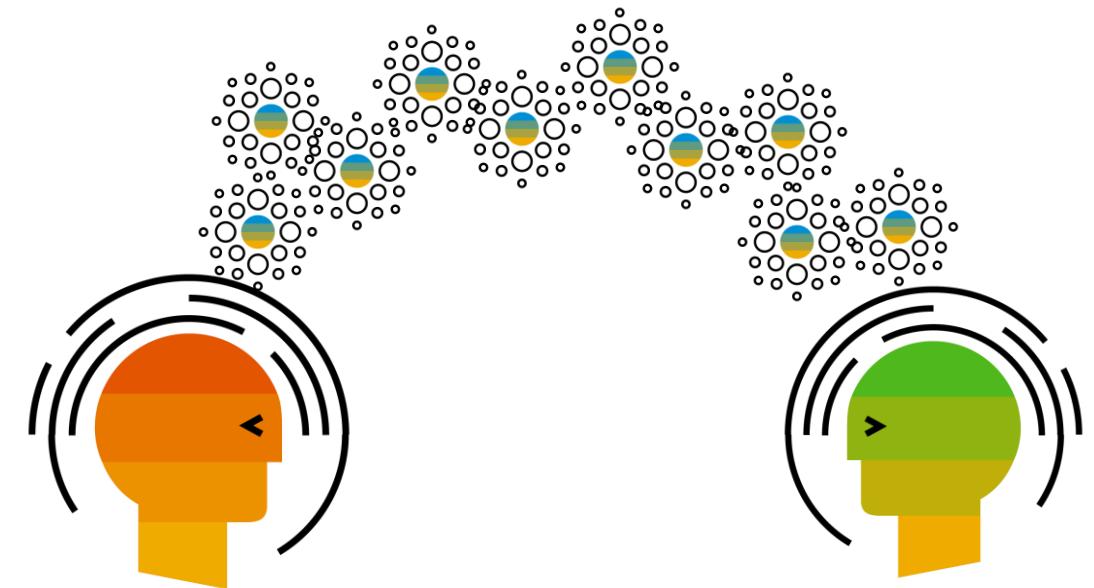


# Accelerating Deep CNN Training

## Overview

### Content:

- Computational considerations
- Batch normalization
- **Transfer learning**
- Residual networks



# Accelerating Deep CNN Training

Transfer learning: motivation

**Problem:** Deep neural networks require very large data sets to train, and large computational resources to build from scratch



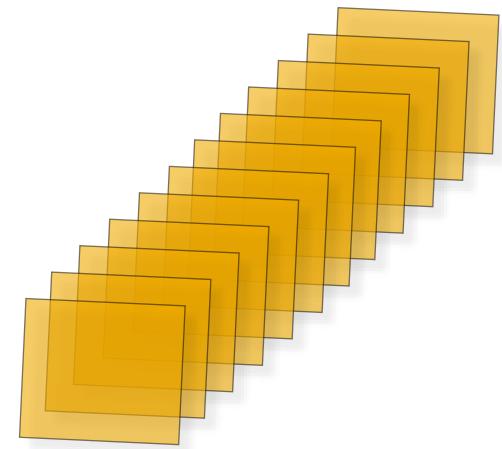
# Accelerating Deep CNN Training

Transfer learning: motivation

**Problem:** Deep neural networks require very large data sets to train, and large computational resources to build from scratch

Famous CNN models **AlexNet**, **ZFNet**, **VGGNet**,  
**GoogleNet**, **Microsoft ResNet...**

- Trained on millions of images
- Powerful GPUs used for training
- Training takes days or even weeks



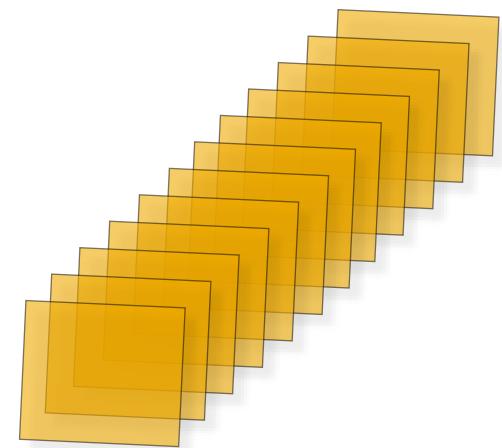
# Accelerating Deep CNN Training

Transfer learning: motivation

**Problem:** Deep neural networks require very large data sets to train, and large computational resources to build from scratch

Famous CNN models **AlexNet**, **ZFNet**, **VGGNet**,  
**GoogleNet**, **Microsoft ResNet...**

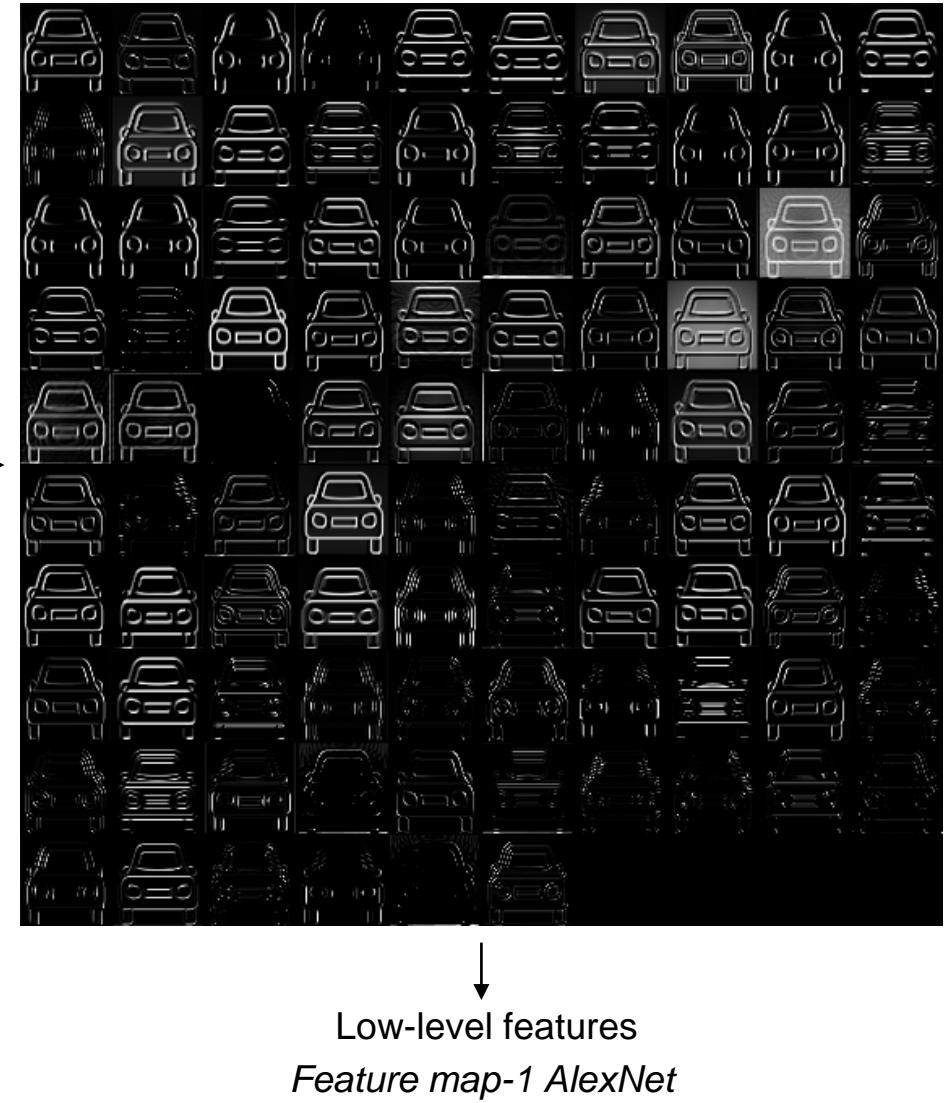
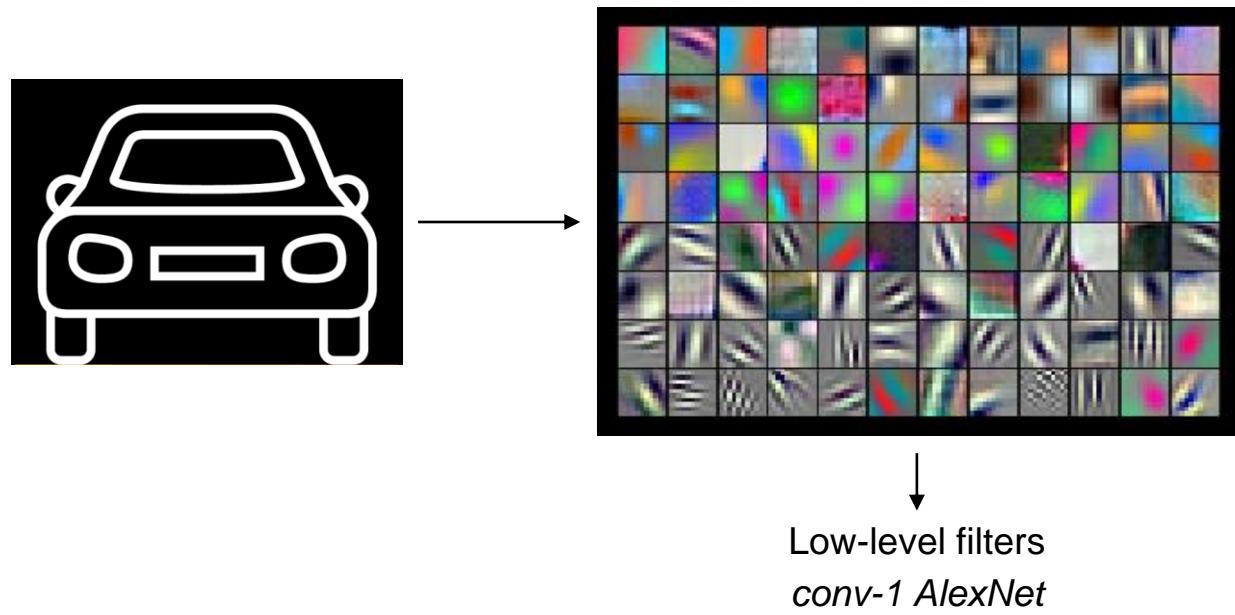
- Trained on millions of images
- Powerful GPUs used for training
- Training takes days or even weeks



**Solution:** Transfer learning (recycle pretrained model)

# Accelerating Deep CNN Training

Transfer learning: reusing joint features

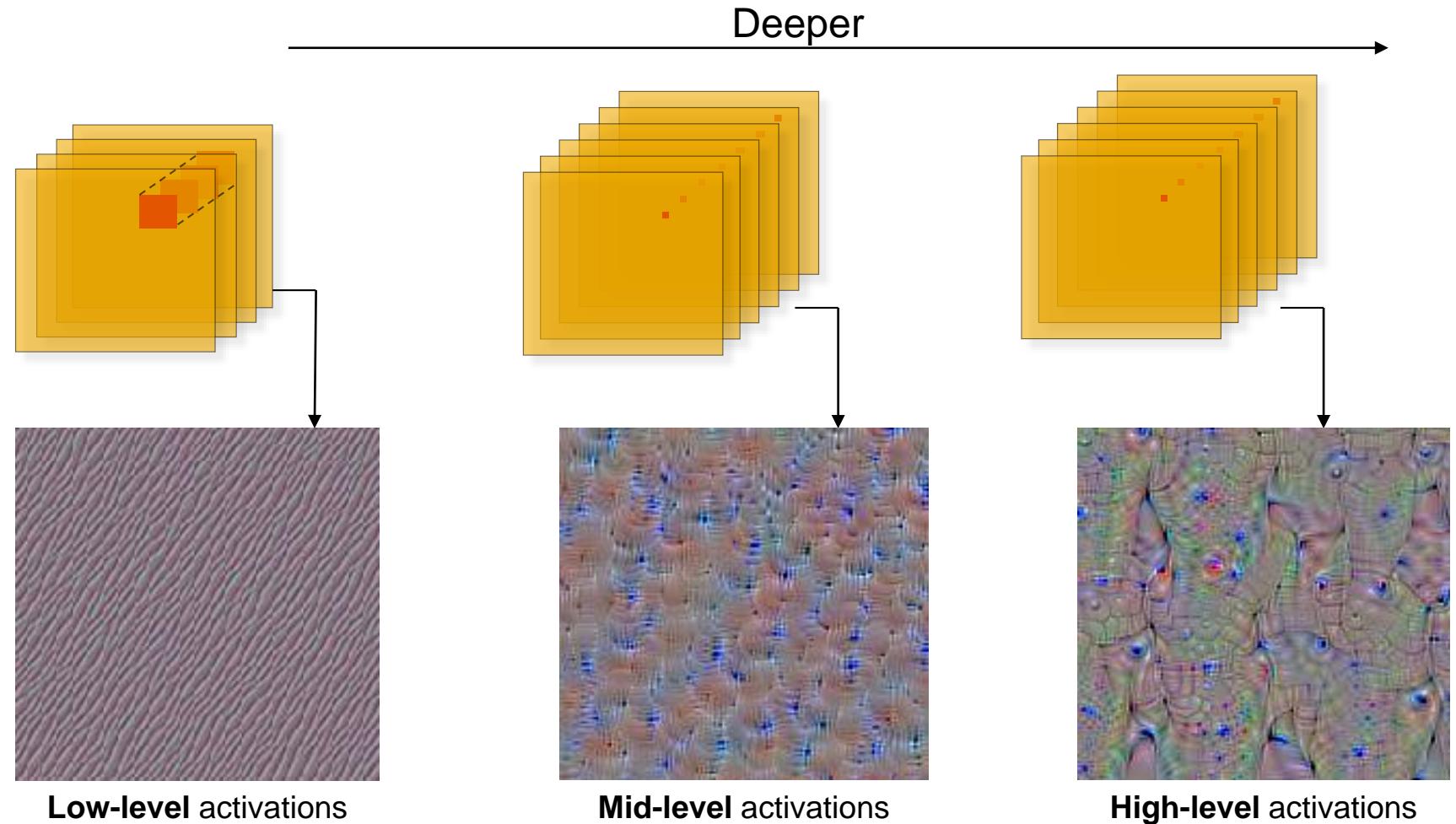


**DeCAF:** Donahue, J. et al. (2014). DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. *CoRR*, abs/1310.1531

**AlexNet:** Krizhevsky, A. et al. (2012). ImageNet Classification with Deep Convolutional Neural Networks. *Commun. ACM* 55, 6.

# Accelerating Deep CNN Training

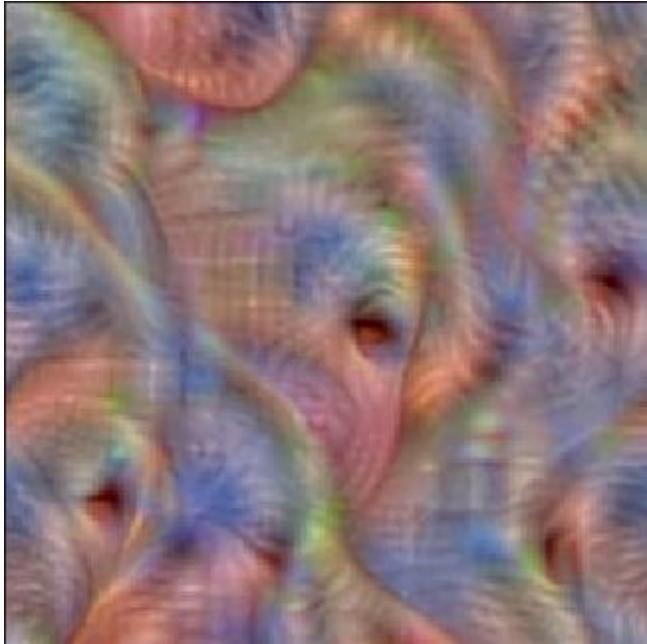
Transfer learning: reusing joint features



VGG Net: Simonyan, K. & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR, abs/1409.1556

# Accelerating Deep CNN Training

Transfer learning: last layer activations



Dog



Cat



Elephant

VGG Net: Simonyan, K. & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR, abs/1409.1556

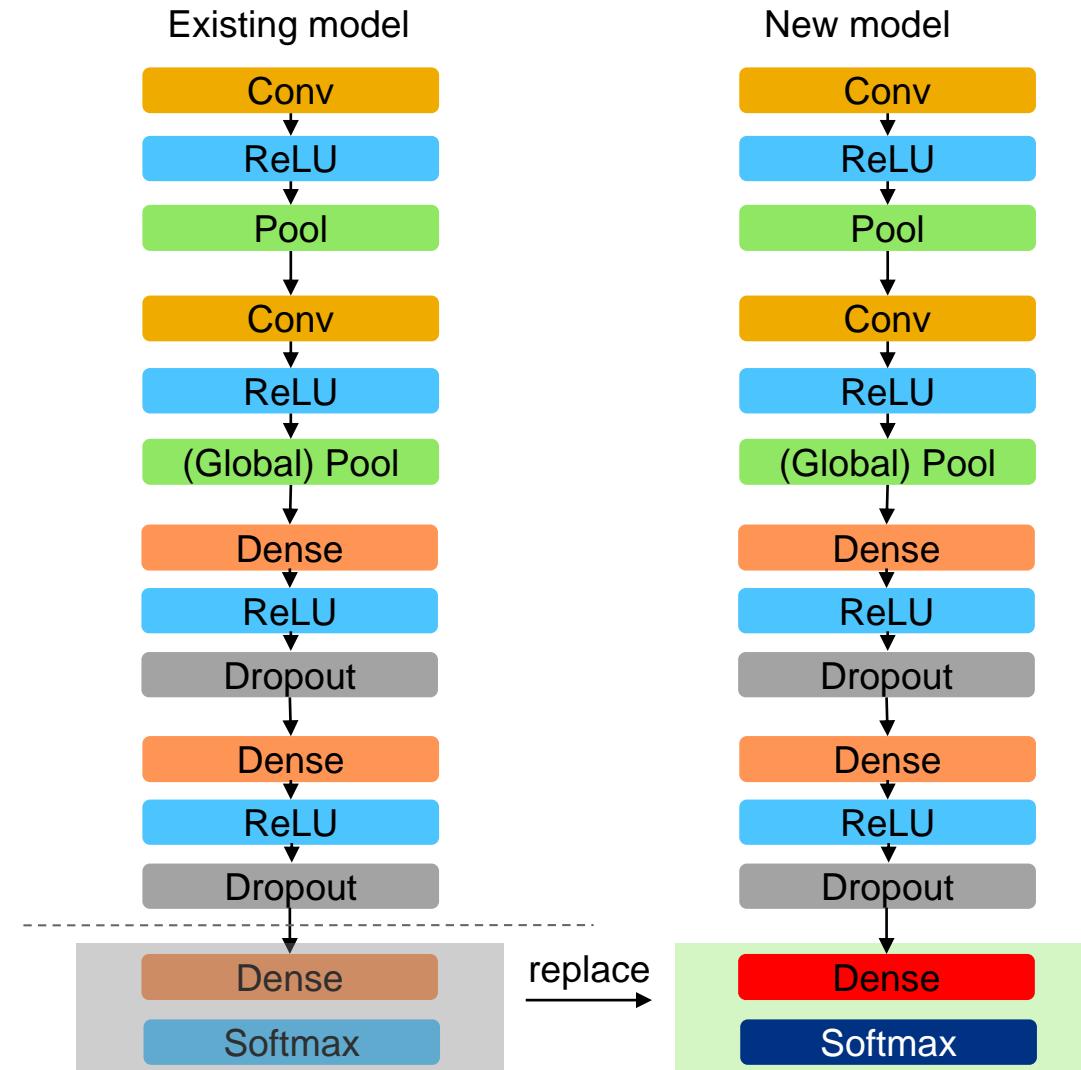
# Accelerating Deep CNN Training

Transfer learning: application

## Scenario 1:

*Replacing last layer with new classifier*

- Retrain on small data set
- Similar features



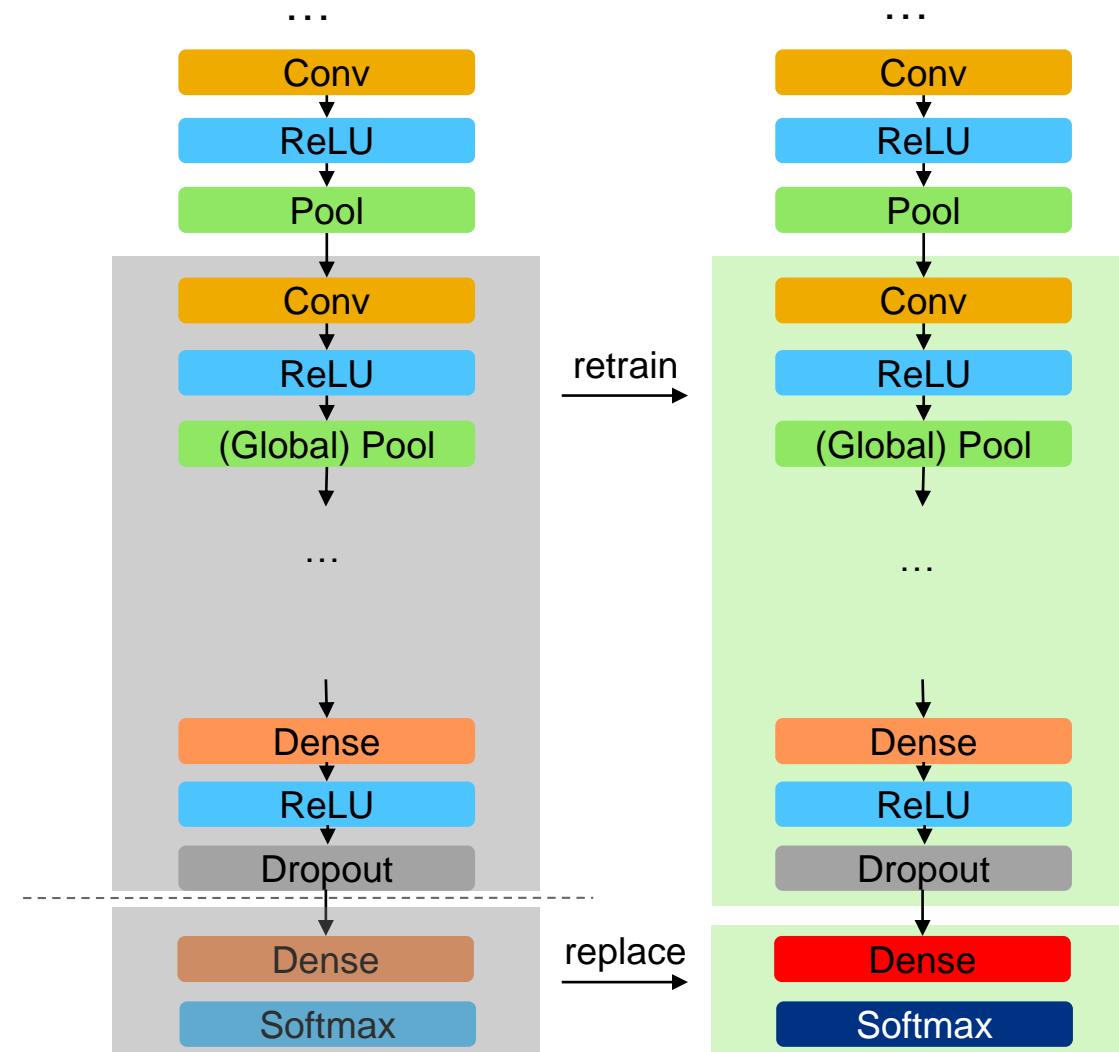
# Accelerating Deep CNN Training

Transfer learning: application

## Scenario 1:

*Fine-tuning more layers*

- Retrain on large data set
- Similar features

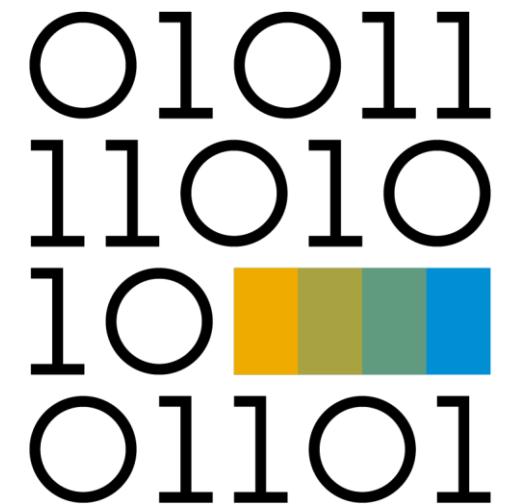


# Accelerating Deep CNN Training

## Overview

### Content:

- Computational considerations
- Batch normalization
- Transfer learning
- **Residual networks**



01011  
11010  
10 █ █ █ █  
01101

# Accelerating Deep CNN Training

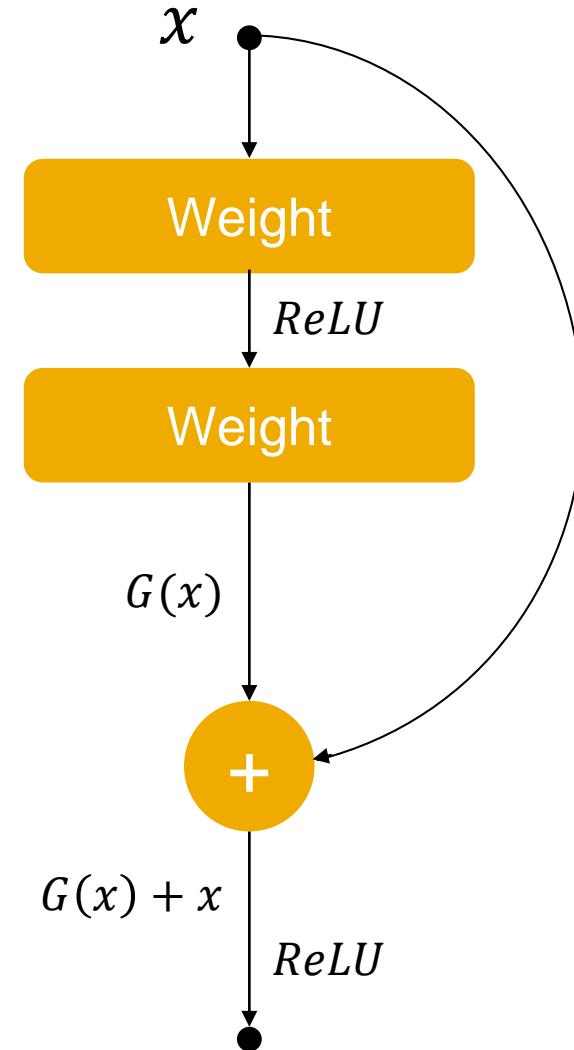
Deep residual networks: Extremely Deep Networks [He, Zhang, Ren, Sun, 2015]

- **Intuition:** Deeper networks are more expressive
- **Problem:** Deep networks are hard to train

# Accelerating Deep CNN Training

Deep residual networks: Extremely Deep Networks [He, Zhang, Ren, Sun, 2015]

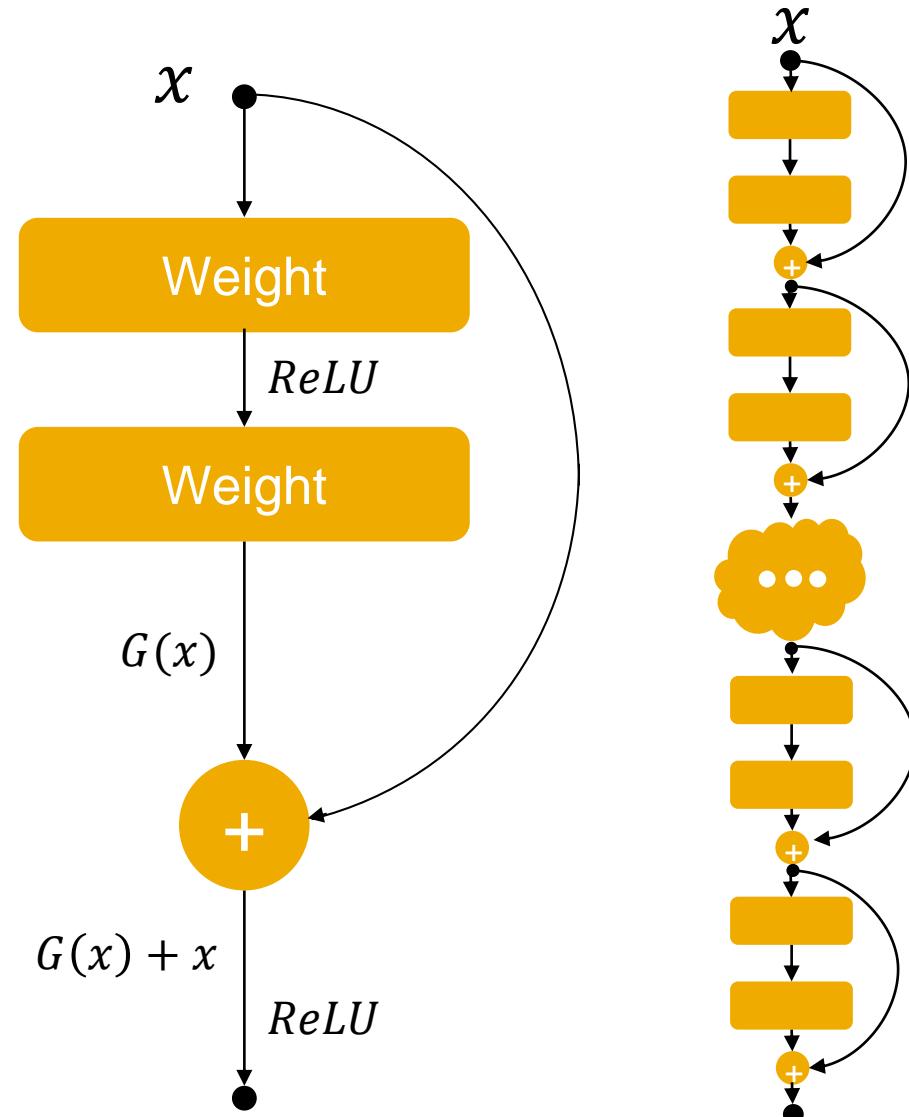
- **Intuition:** Deeper networks are more expressive
- **Problem:** Deep networks are hard to train
- **Idea:** The modeled function of each building block has a higher resemblance to the identity function (than to the zero function)



# Accelerating Deep CNN Training

Deep residual networks: Extremely Deep Networks [He, Zhang, Ren, Sun, 2015]

- **Intuition:** Deeper networks are more expressive
- **Problem:** Deep networks are hard to train
- **Idea:** The modeled function of each building block has a higher resemblance to the identity function (than to the zero function)
- The signal can *directly skip* through multiple layers (gradient more easily flows from top to bottom layers)
- Allows very deep architectures (typically more than 100 layers)

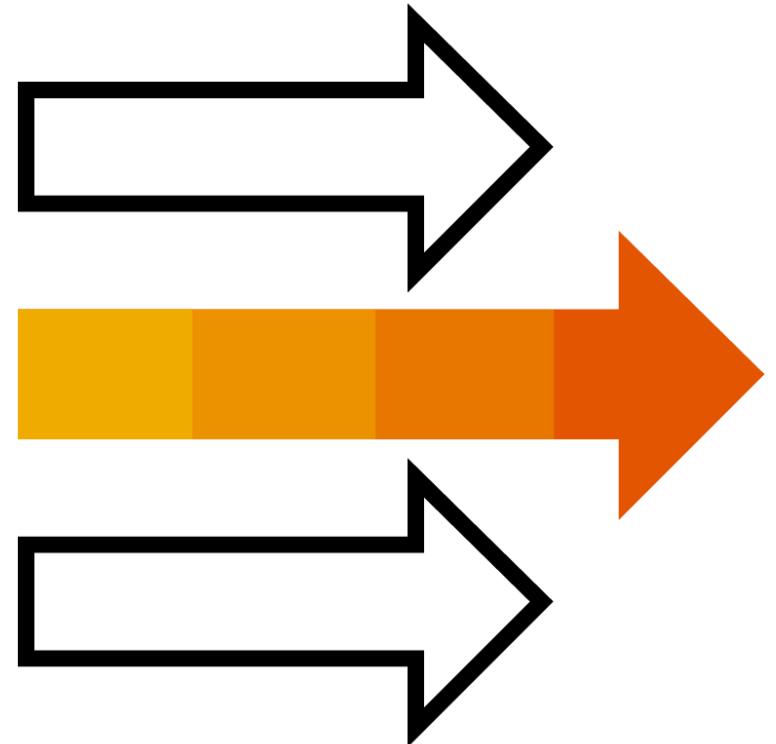


# Accelerating Deep CNN Training

Coming up next

## Applications of CNNs

- Object detection
- Semantic image segmentation



# Thank you.

Contact information:

[open@sap.com](mailto:open@sap.com)

# © 2017 SAP SE or an SAP affiliate company. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP SE or an SAP affiliate company.

The information contained herein may be changed without prior notice. Some software products marketed by SAP SE and its distributors contain proprietary software components of other software vendors. National product specifications may vary.

These materials are provided by SAP SE or an SAP affiliate company for informational purposes only, without representation or warranty of any kind, and SAP or its affiliated companies shall not be liable for errors or omissions with respect to the materials. The only warranties for SAP or SAP affiliate company products and services are those that are set forth in the express warranty statements accompanying such products and services, if any. Nothing herein should be construed as constituting an additional warranty.

In particular, SAP SE or its affiliated companies have no obligation to pursue any course of business outlined in this document or any related presentation, or to develop or release any functionality mentioned therein. This document, or any related presentation, and SAP SE's or its affiliated companies' strategy and possible future developments, products, and/or platform directions and functionality are all subject to change and may be changed by SAP SE or its affiliated companies at any time for any reason without notice. The information in this document is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. All forward-looking statements are subject to various risks and uncertainties that could cause actual results to differ materially from expectations. Readers are cautioned not to place undue reliance on these forward-looking statements, and they should not be relied upon in making purchasing decisions.

SAP and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP SE (or an SAP affiliate company) in Germany and other countries. All other product and service names mentioned are the trademarks of their respective companies.

See <http://global.sap.com/corporate-en/legal/copyright/index.epx> for additional trademark information and notices.



Week 4: Convolutional Networks

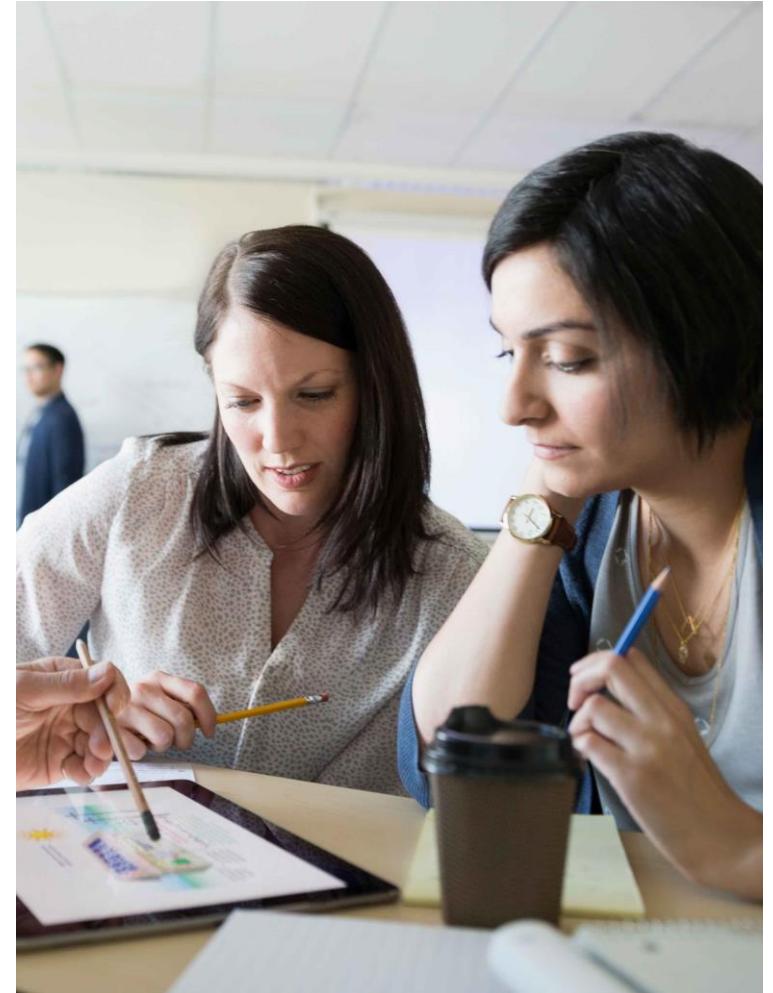
## **Unit 5: Applications of CNNs**

# Applications of CNNs

What we covered in the last unit

## Accelerating Deep CNN Training

- Computational considerations
- Batch normalization
- Transfer learning
- Residual networks

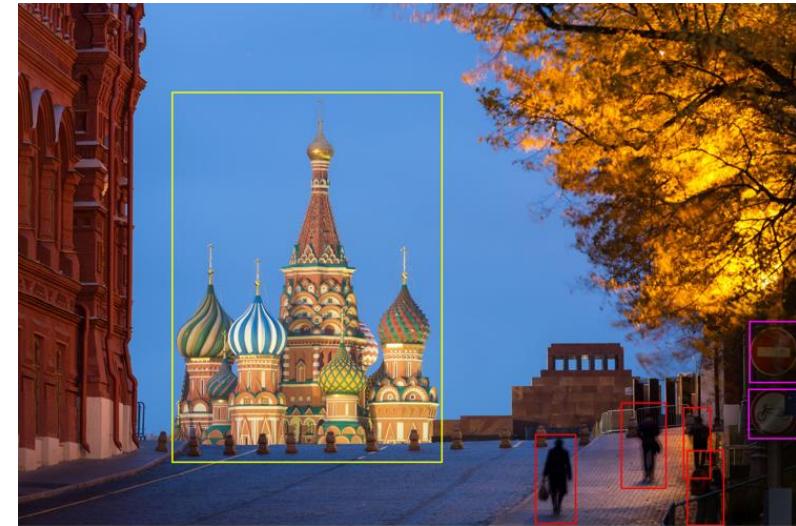


# Applications of CNNs

## Overview

### Content:

- Object detection
  - Two-stage detectors
  - One-stage detectors
- Segmentation

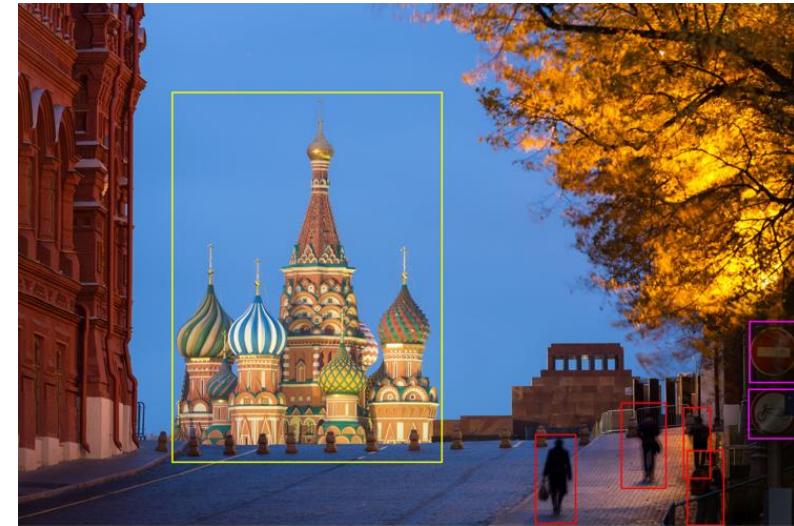


# Applications of CNNs

## Overview

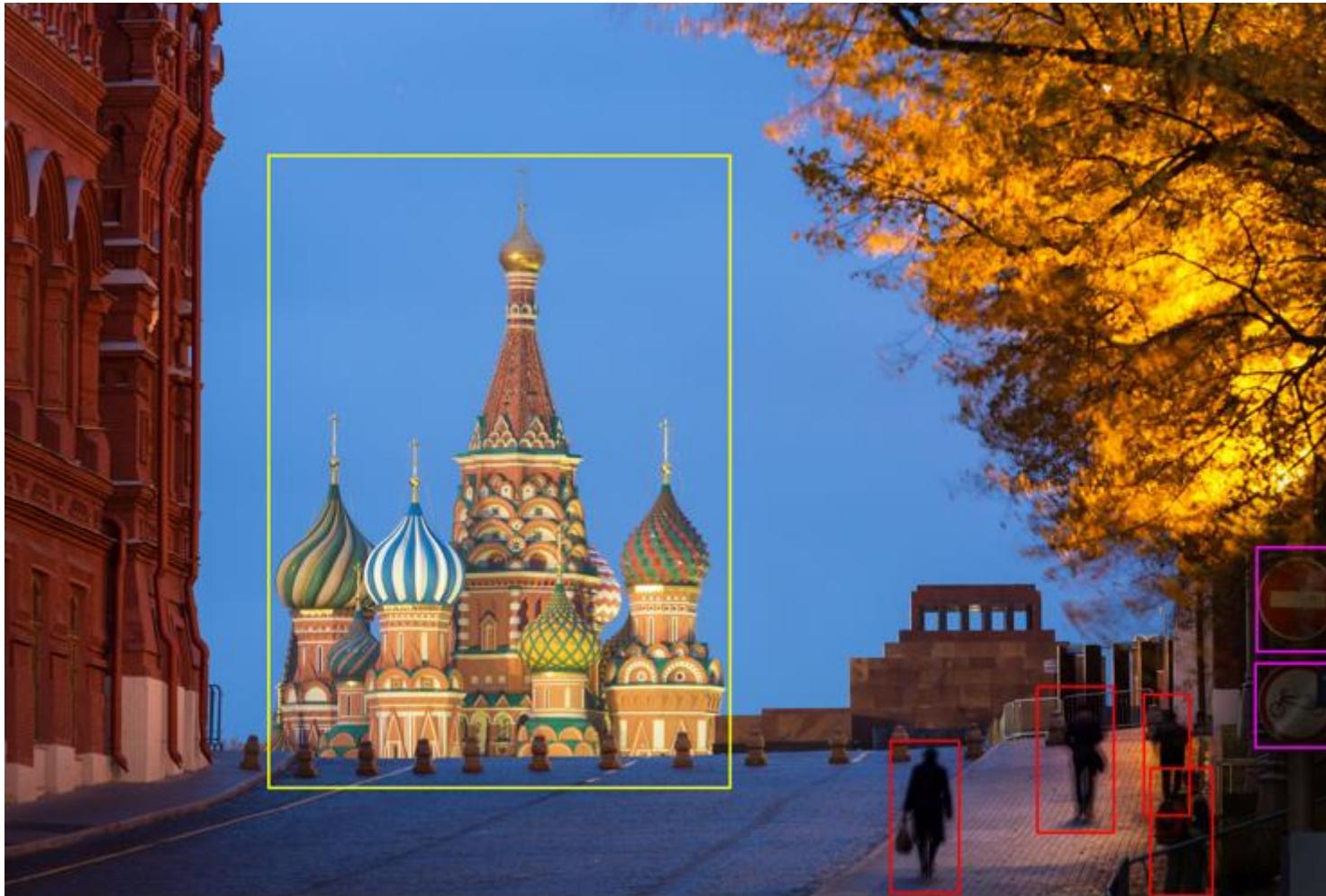
### Content:

- Object detection
  - Two-stage detectors
  - One-stage detectors
- Segmentation



# Applications of CNNs

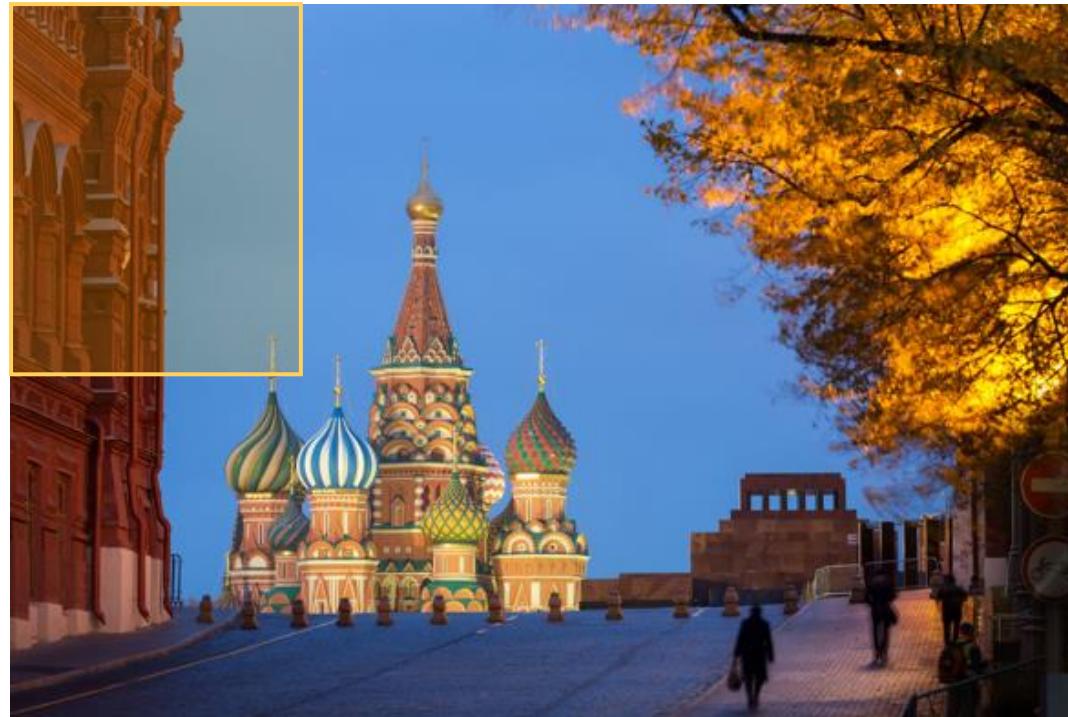
Object detection: example



# Applications of CNNs

Object detection – Naïve approach

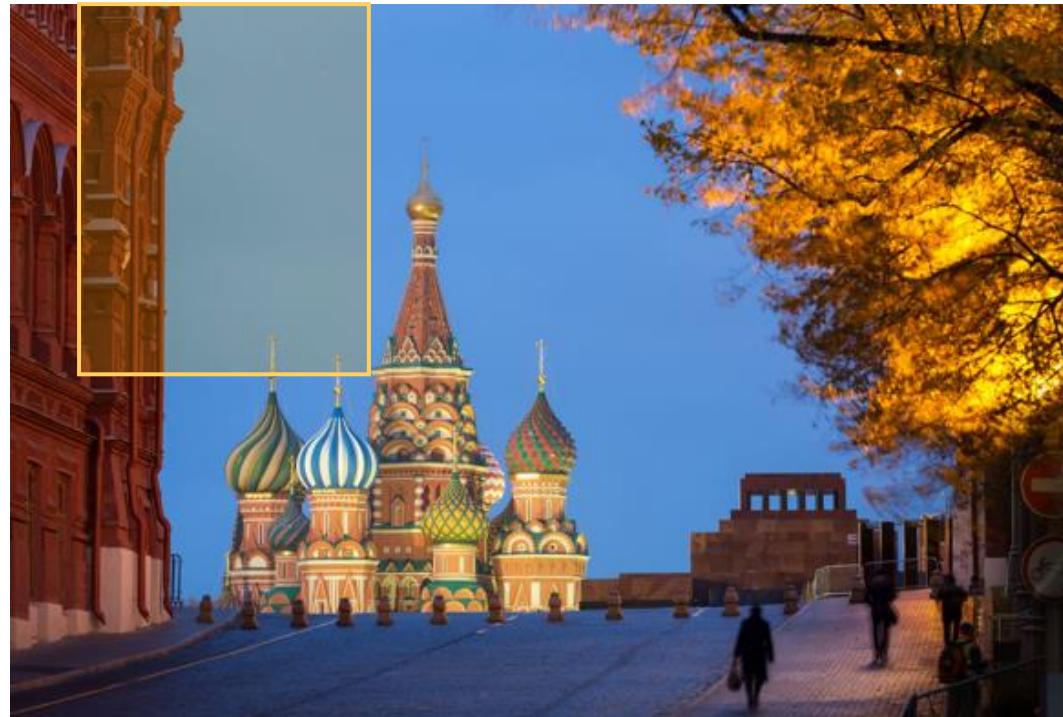
## 1. Specify sliding window size



# Applications of CNNs

Object detection – Naïve approach

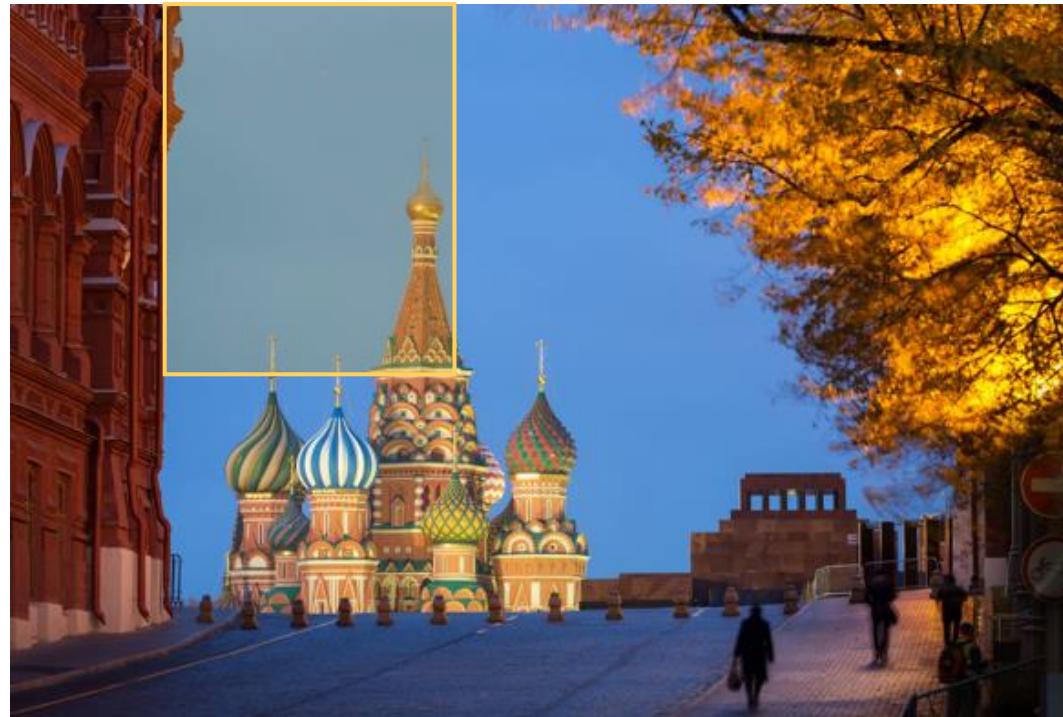
## 1. Specify sliding window size



# Applications of CNNs

Object detection – Naïve approach

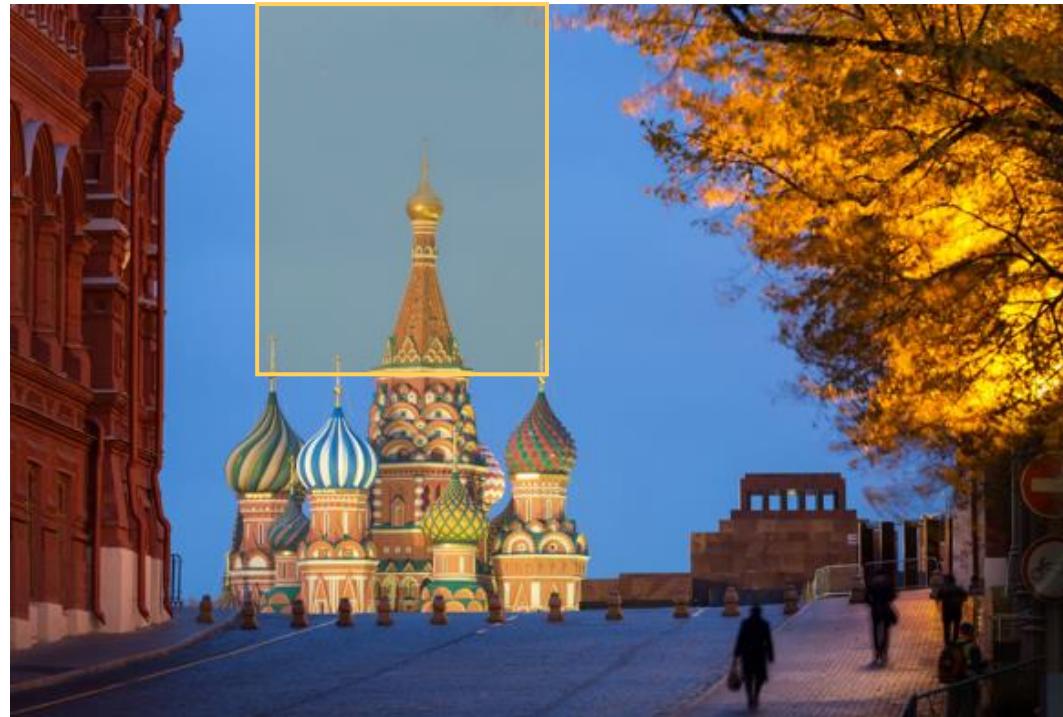
## 1. Specify sliding window size



# Applications of CNNs

Object detection – Naïve approach

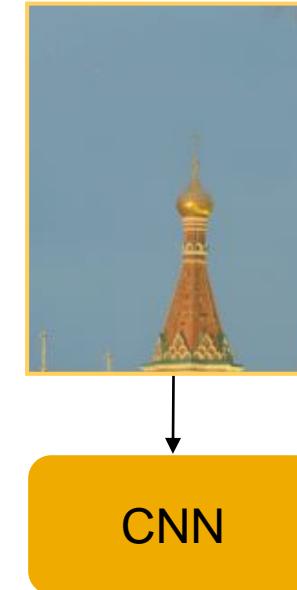
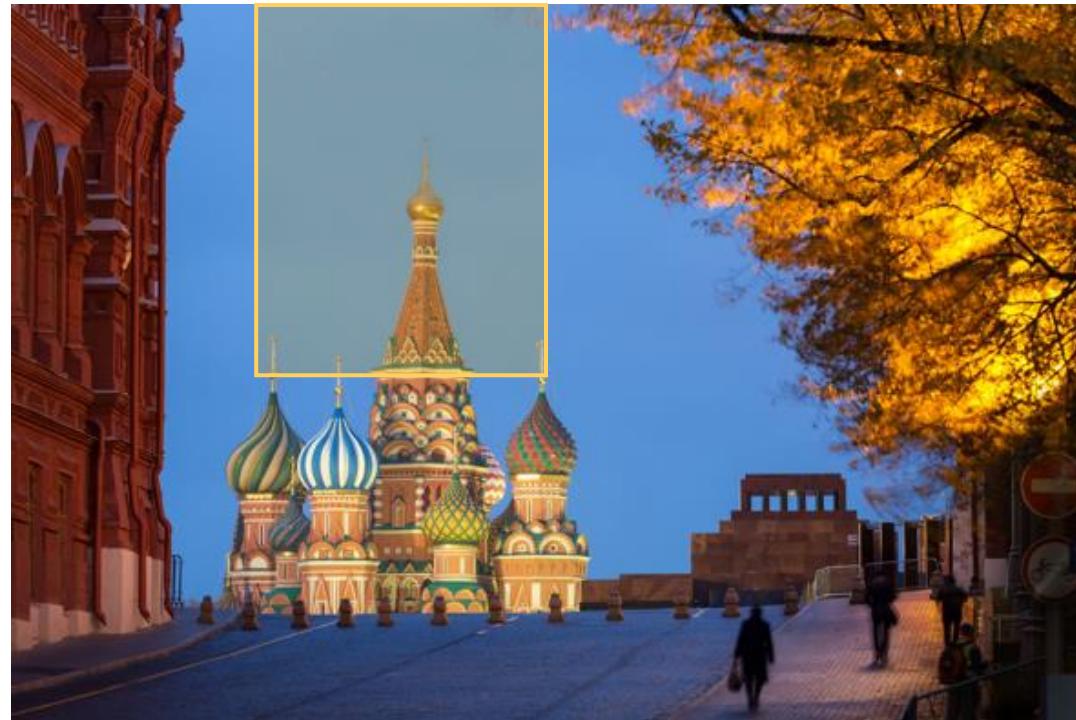
## 1. Specify sliding window size



# Applications of CNNs

Object detection – Naïve approach

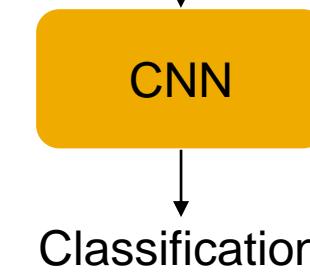
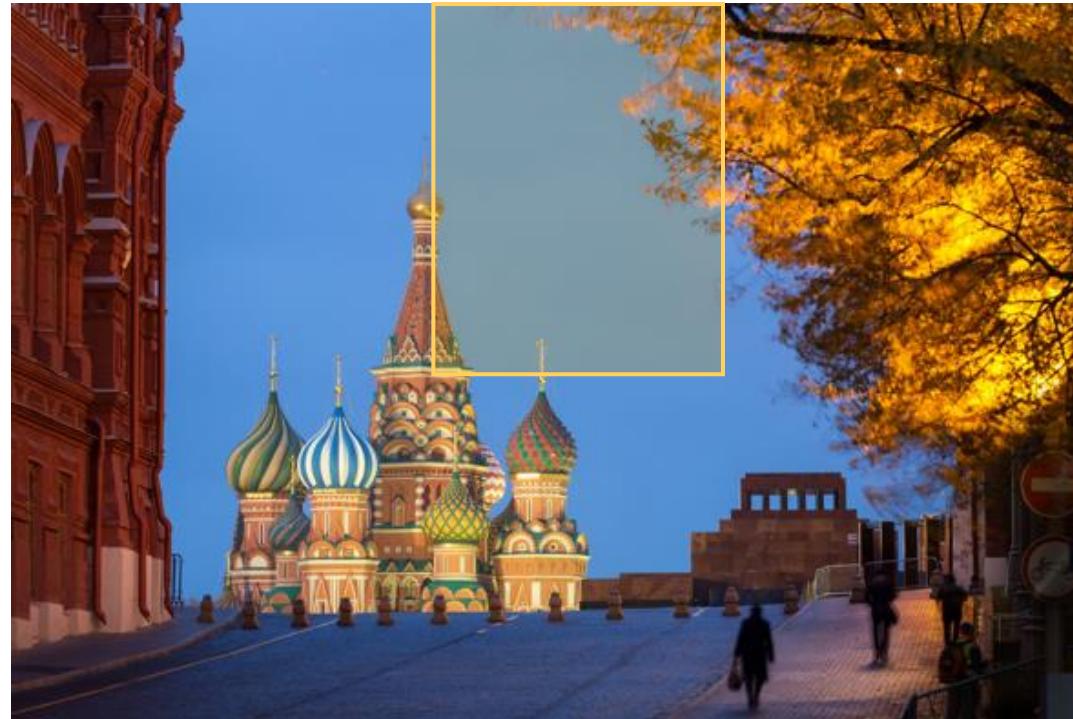
1. Specify sliding window size
2. Run pre-trained object-classification CNN on each window (*softmax* output)



# Applications of CNNs

## Object detection – Naïve approach

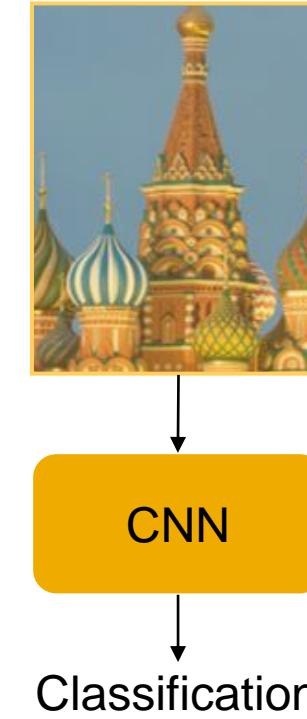
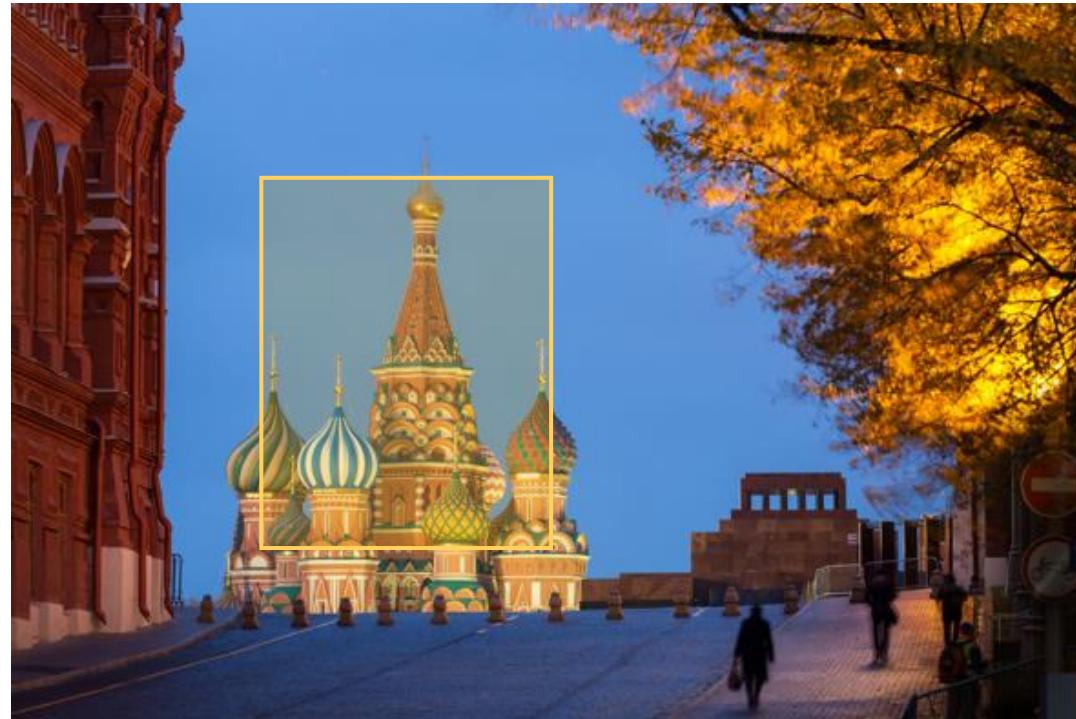
1. Specify sliding window size
2. Run pre-trained object-classification CNN on each window (*softmax* output)
3. Associate current sliding window with object class if probability is sufficiently high (e.g. > 0.5)



# Applications of CNNs

Object detection – Naïve approach

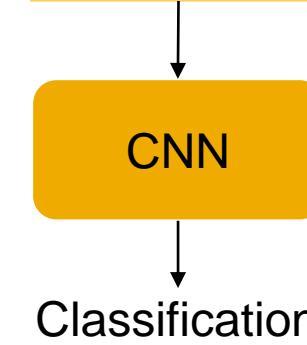
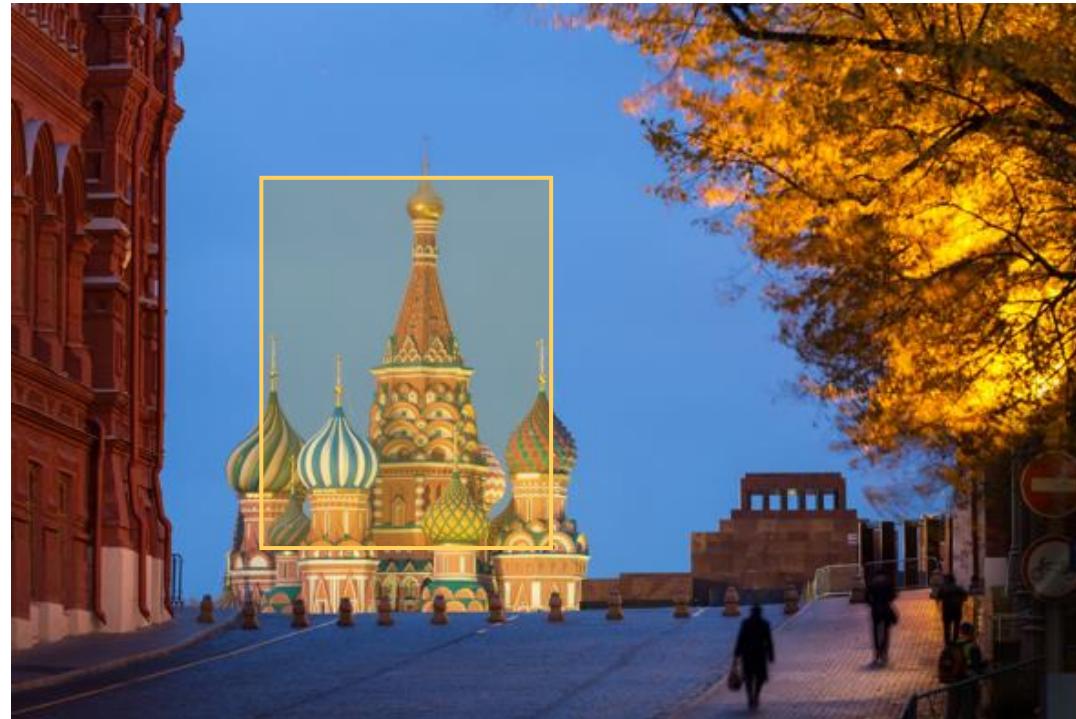
1. Specify sliding window size
2. Run pre-trained object-classification CNN on each window (*softmax* output)
3. Associate current sliding window with object class if probability is sufficiently high (e.g. > 0.5)



# Applications of CNNs

Object detection – Naïve approach

1. Specify sliding window size
2. Run pre-trained object-classification CNN on each window (softmax output)
3. Associate current sliding window with object class if probability is sufficiently high (e.g.  $> 0.5$ )



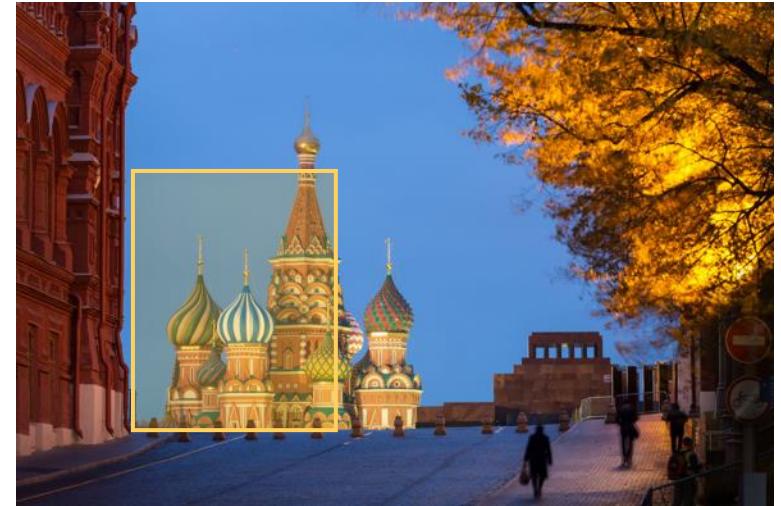
## Problem:

*Intractable even  
for small images*

# Applications of CNNs

Two-stage detectors: key idea

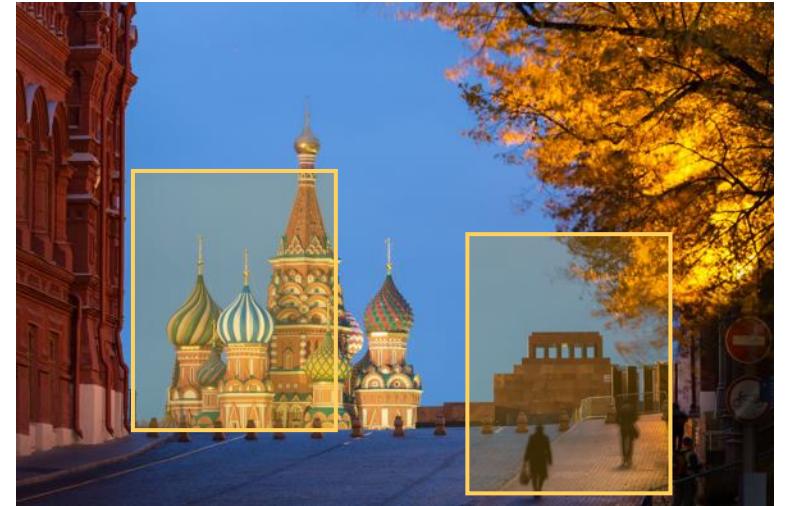
- Use proposal algorithm (e.g. "selective search")



# Applications of CNNs

Two-stage detectors: key idea

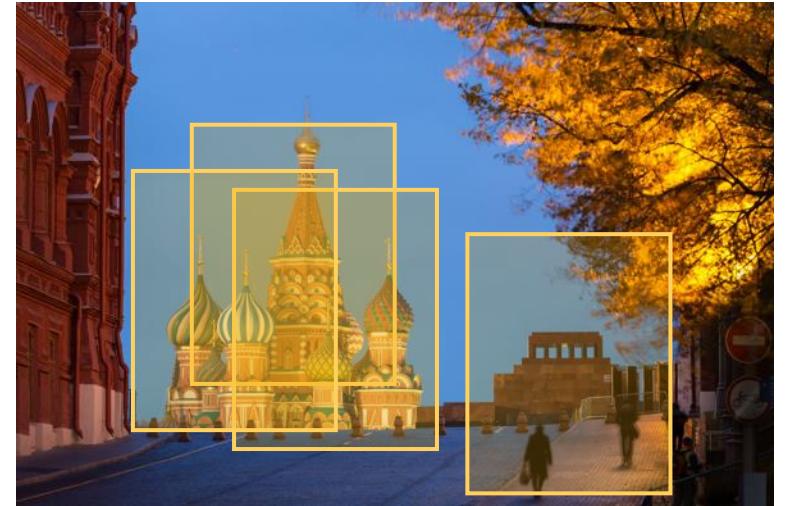
- Use proposal algorithm (e.g. “selective search”)
- Boxes are proposed based on e.g.:
  - Texture
  - Color
  - Intensity



# Applications of CNNs

Two-stage detectors: key idea

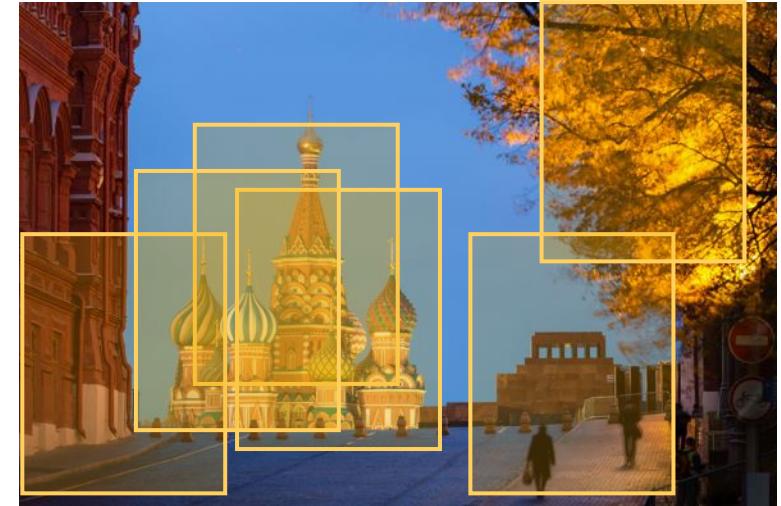
- Use proposal algorithm (e.g. “selective search”)
- Boxes are proposed based on e.g.:
  - Texture
  - Color
  - Intensity



# Applications of CNNs

Two-stage detectors: key idea

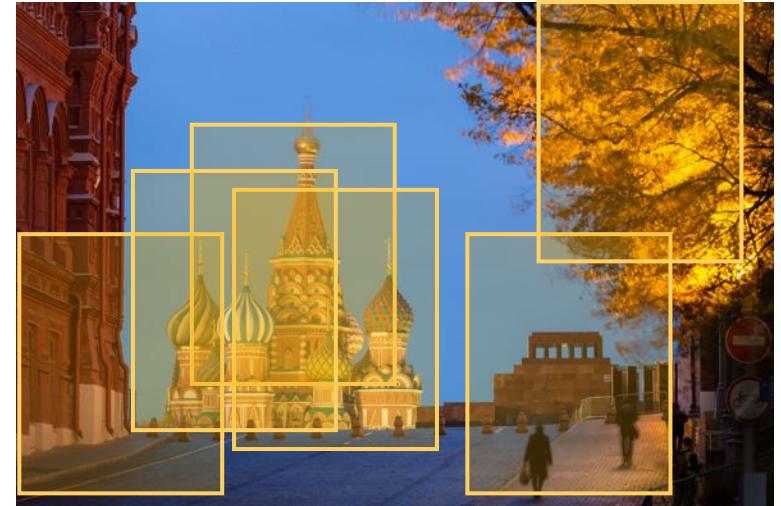
- Use proposal algorithm (e.g. “selective search”)
- Boxes are proposed based on e.g.:
  - Texture
  - Color
  - Intensity



# Applications of CNNs

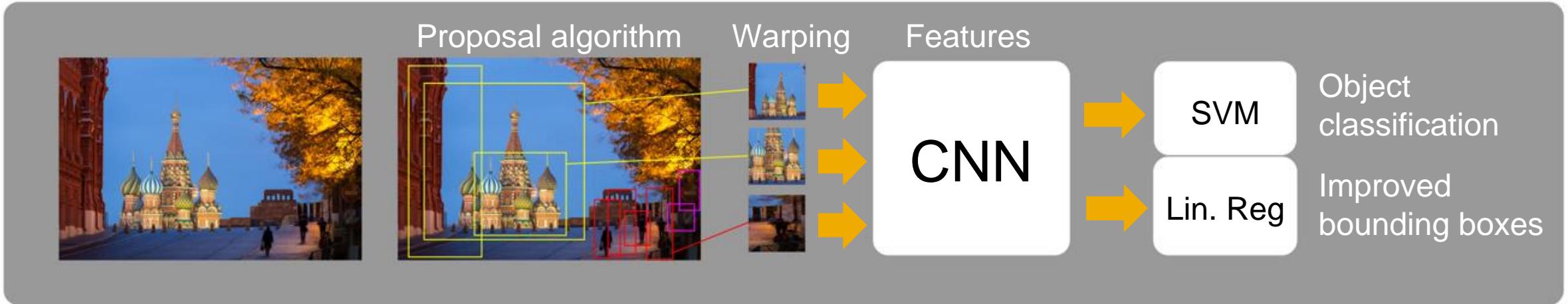
Two-stage detectors: key idea

- Use proposal algorithm (e.g. “selective search”)
- Boxes are proposed based on e.g.:
  - Texture
  - Color
  - Intensity
- Evaluate only on “proposals”



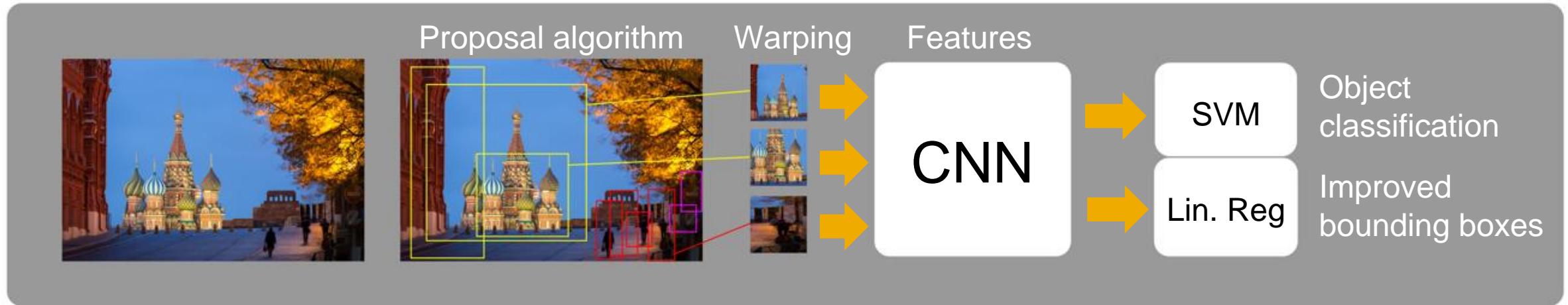
# Applications of CNNs

Two-stage detectors: R-CNN [Girshick, Donahue, Darrell, Malik, 2013]



# Applications of CNNs

Two-stage detectors: R-CNN [Girshick, Donahue, Darrell, Malik, 2013]

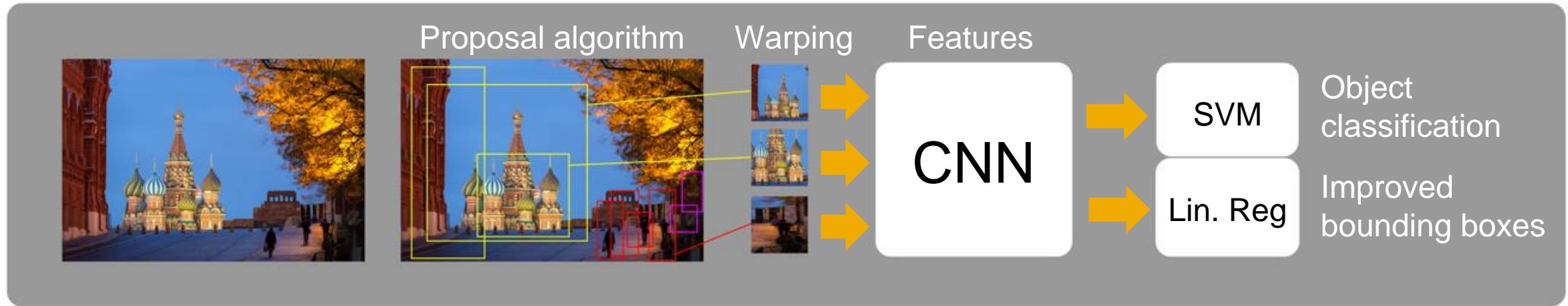


## Advantages:

- Several orders of magnitude faster
- Allows improvement of bounding boxes

# Applications of CNNs

Two-stage detectors: R-CNN [Girshick, Donahue, Darrell, Malik, 2013]



## Advantages:

- Several orders of magnitude faster
- Allows improvement of bounding boxes

## Disadvantages:

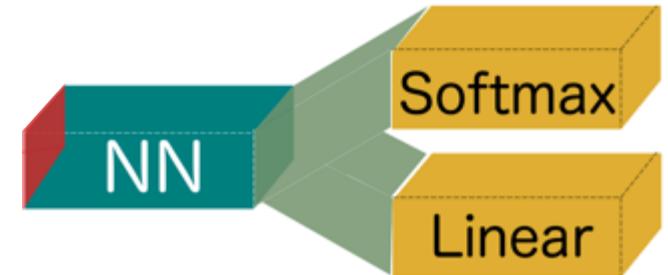
- CNN cannot adapt to changes in SVM/linear regression
  - *Implication:* Degrades power of model
- CNN often evaluates boxes that have large overlaps
  - *Implication:* Suboptimal performance

# Applications of CNNs

Two-stage detectors: Fast R-CNN [Girshick 2015]

## Key Idea:

- Softmax output replaces SVM
- Combine regression and classification into one neural net (end-to-end)



Fully connected

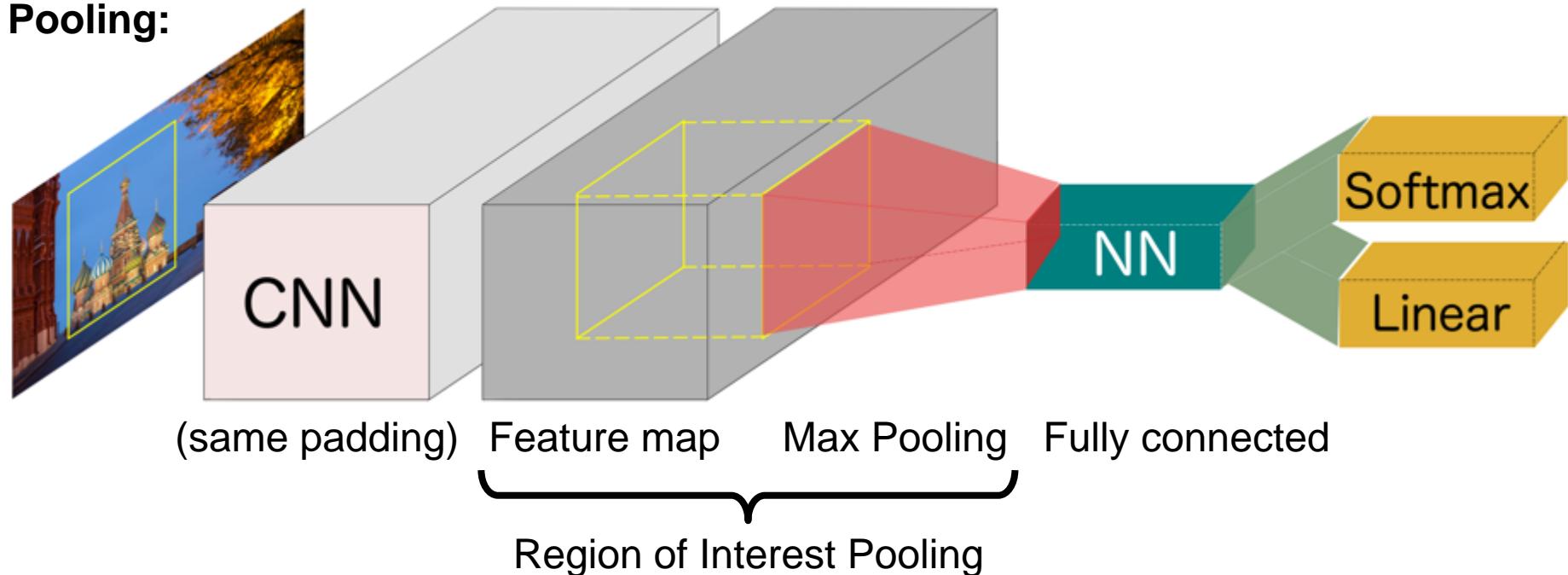
# Applications of CNNs

Two-stage detectors: Fast R-CNN [Girshick 2015]

## Key Idea:

- Softmax output replaces SVM
- Combine regression and classification into one neural net (end-to-end)
- Evaluate CNN only one time/image (share results among suggestions)

## Region of Interest Pooling:



# Applications of CNNs

Two-stage detectors: Faster R-CNN [*Ren, He, Girshick, Sun, 2015*]

**Key Problem:** Region proposal becomes bottleneck

# Applications of CNNs

Two-stage detectors: Faster R-CNN [*Ren, He, Girshick, Sun, 2015*]

**Key Problem:** Region proposal becomes bottleneck

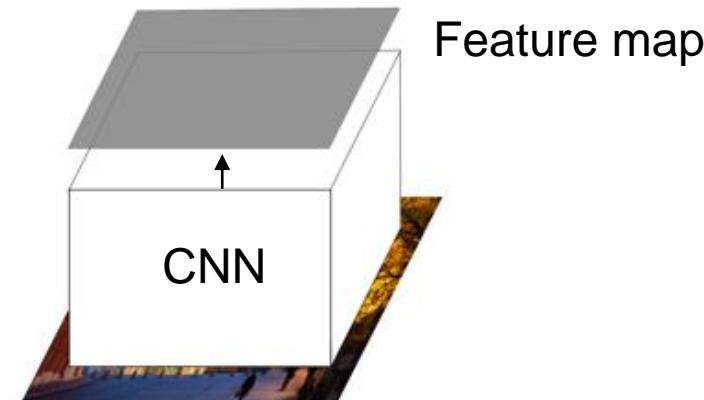
**Solution:** Add so-called “Region Proposal Network”

# Applications of CNNs

Two-stage detectors: Faster R-CNN [Ren, He, Girshick, Sun, 2015]

**Key Problem:** Region proposal becomes bottleneck

**Solution:** Add so-called “Region Proposal Network”

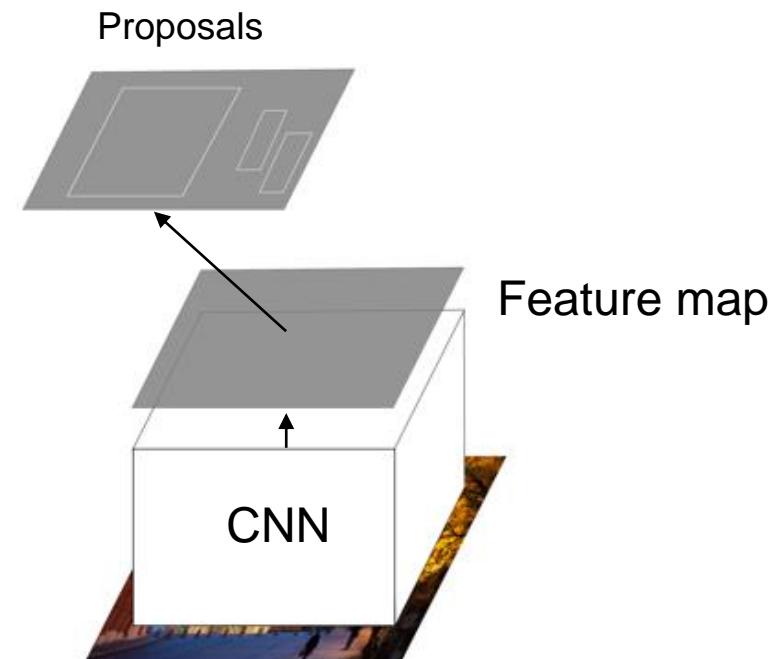


# Applications of CNNs

Two-stage detectors: Faster R-CNN [Ren, He, Girshick, Sun, 2015]

**Key Problem:** Region proposal becomes bottleneck

**Solution:** Add so-called “Region Proposal Network”

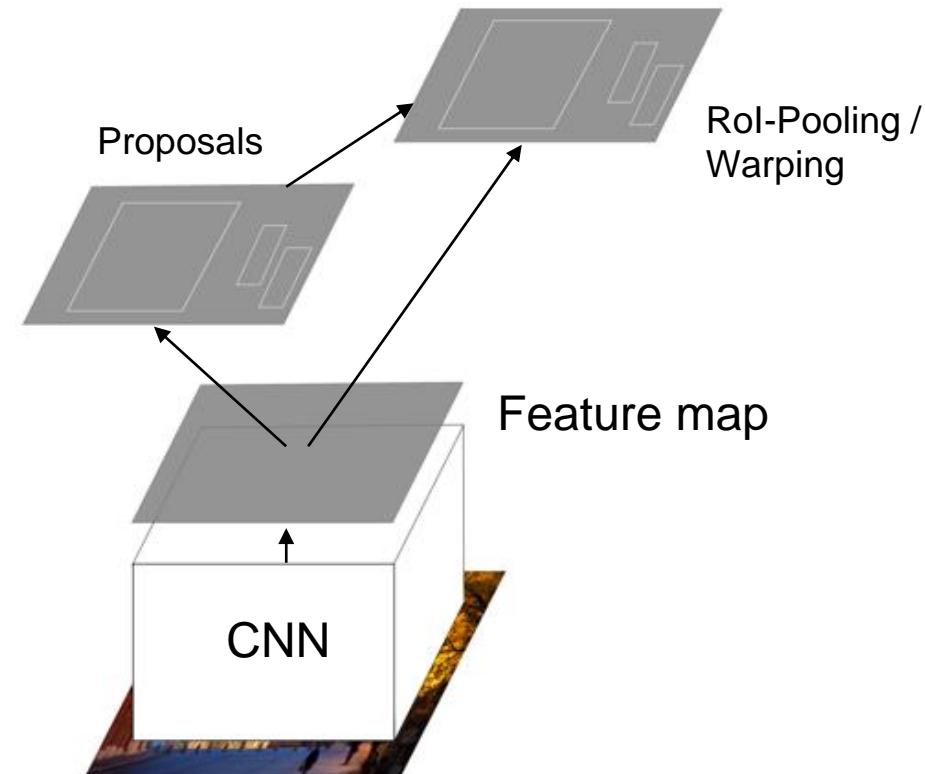


# Applications of CNNs

Two-stage detectors: Faster R-CNN [Ren, He, Girshick, Sun, 2015]

**Key Problem:** Region proposal becomes bottleneck

**Solution:** Add so-called “Region Proposal Network”

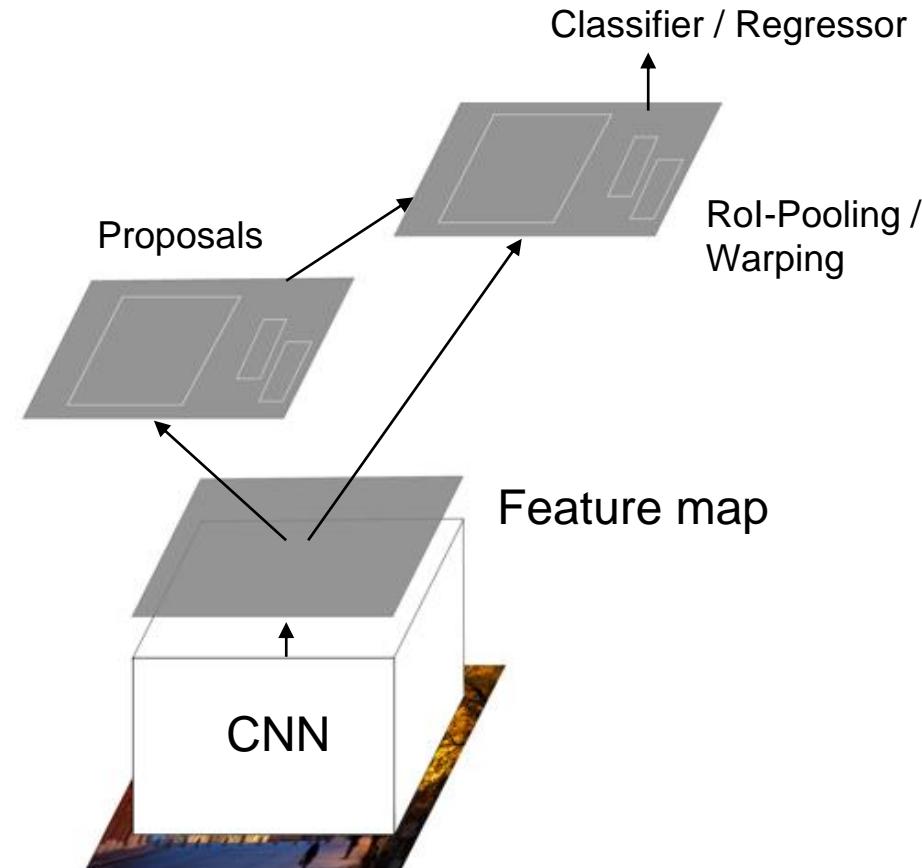


# Applications of CNNs

Two-stage detectors: Faster R-CNN [Ren, He, Girshick, Sun, 2015]

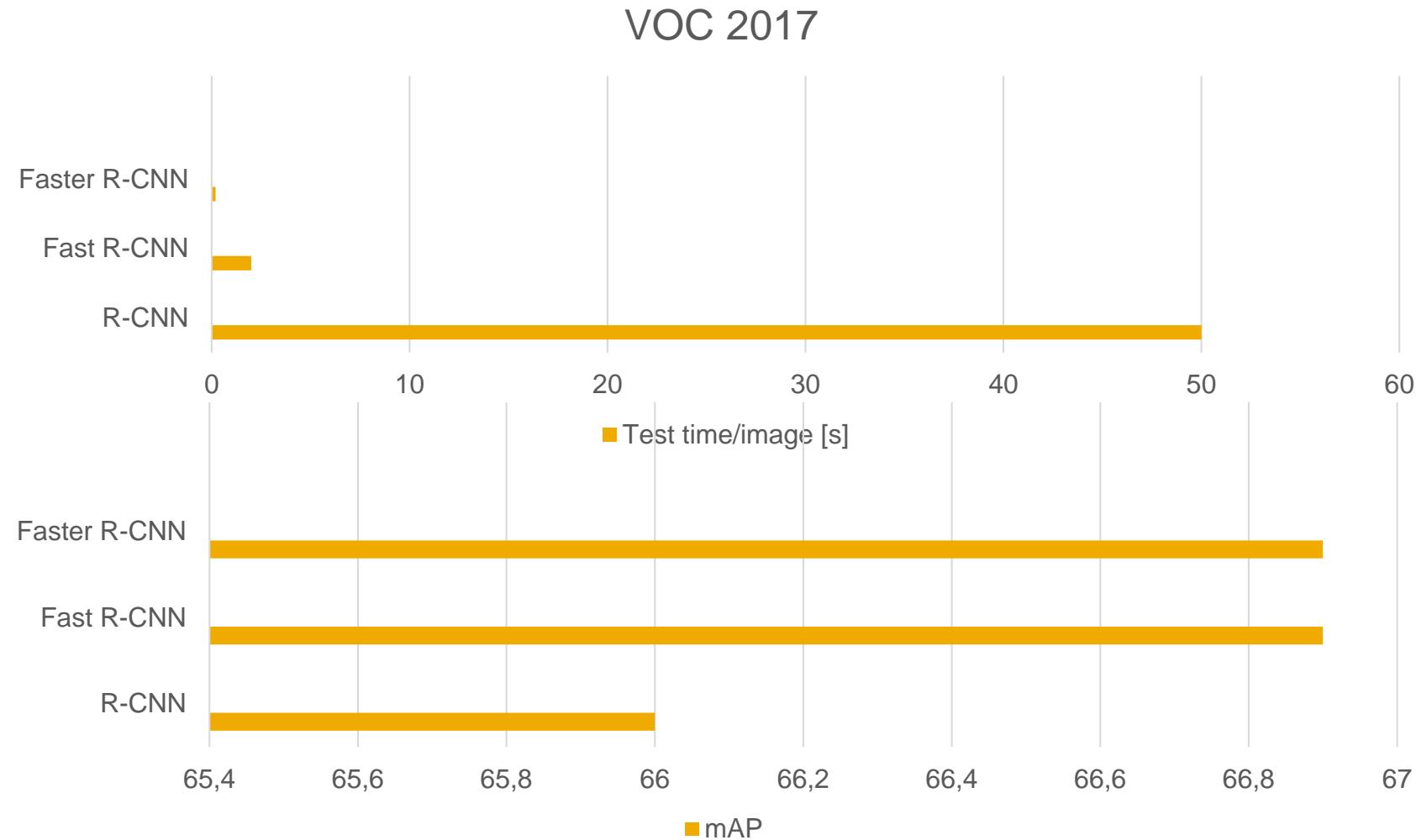
**Key Problem:** Region proposal becomes bottleneck

**Solution:** Add so-called “Region Proposal Network”



# Applications of CNNs

Two-stage detectors: comparison [Li, Karpathy, Johnson]

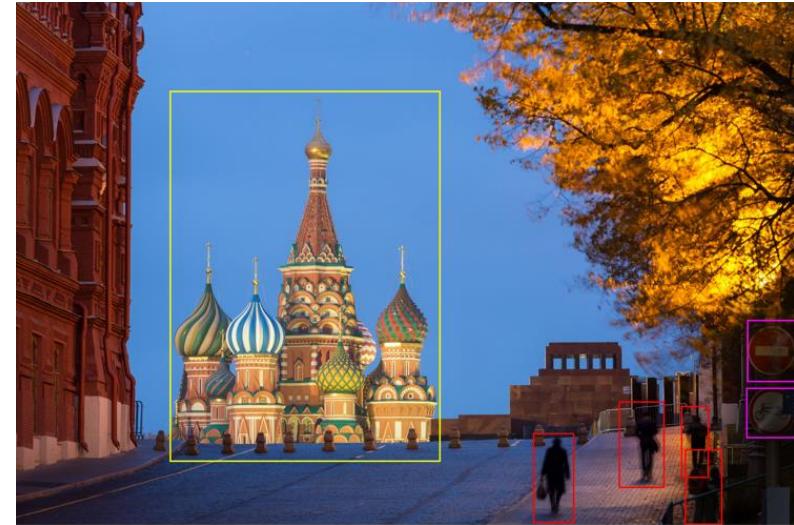


# Applications of CNNs

## Overview

### Content:

- **Object detection**
  - Two-stage detectors
  - **One-stage detectors**
- Segmentation



# Applications of CNNs

Object detection: one-stage detectors

How can we make  
this faster?



The R-CNN class detectors are accurate – but very slow!

**Methodology:**

- First stage: Generate region proposals
- Second stage: Classify & refine bbox proposals

# Applications of CNNs

Object detection: one-stage detectors

How can we make  
this faster?

The R-CNN class detectors are accurate – but very slow!

**Methodology:**

- First stage: Generate region proposals
- Second stage: Classify & refine bbox proposals

CHANGE  
METHODOLOGY

One-stage detectors:

- **Idea:** Take advantage of fully convolutional networks to compute bounding box regression and class probabilities in **one** pass
- **Implication:** Architecture is considerably simpler!

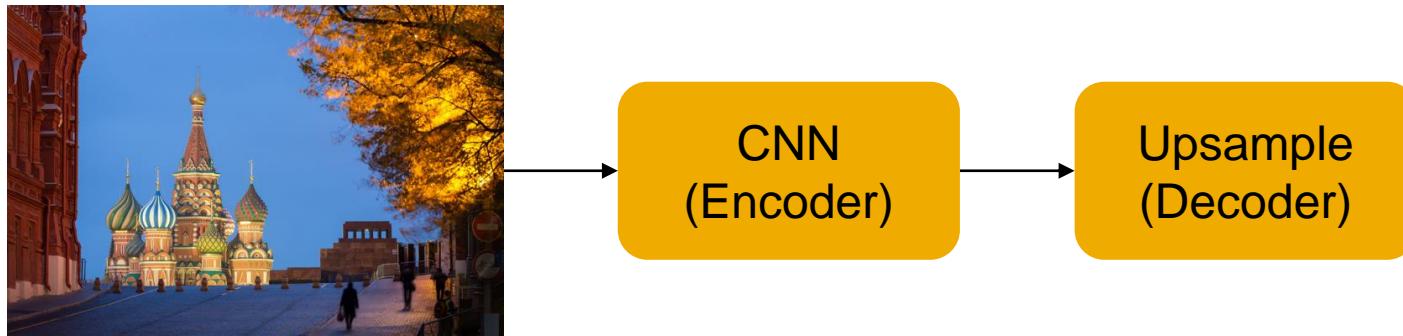
# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]



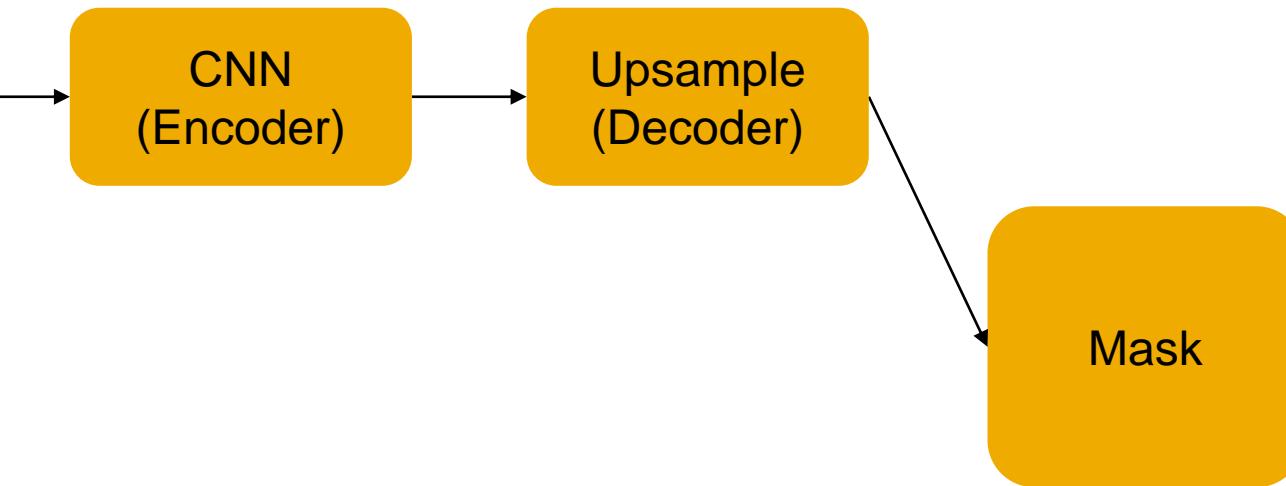
# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]



# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]

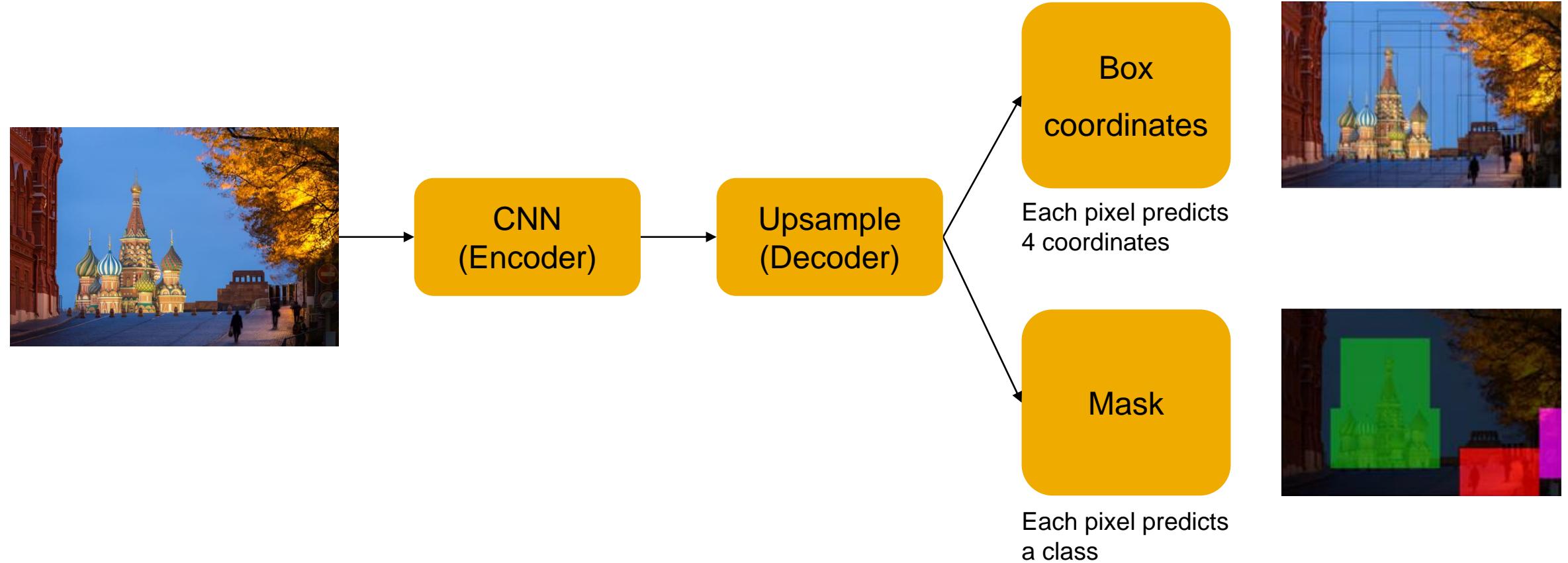


Each pixel predicts  
a class



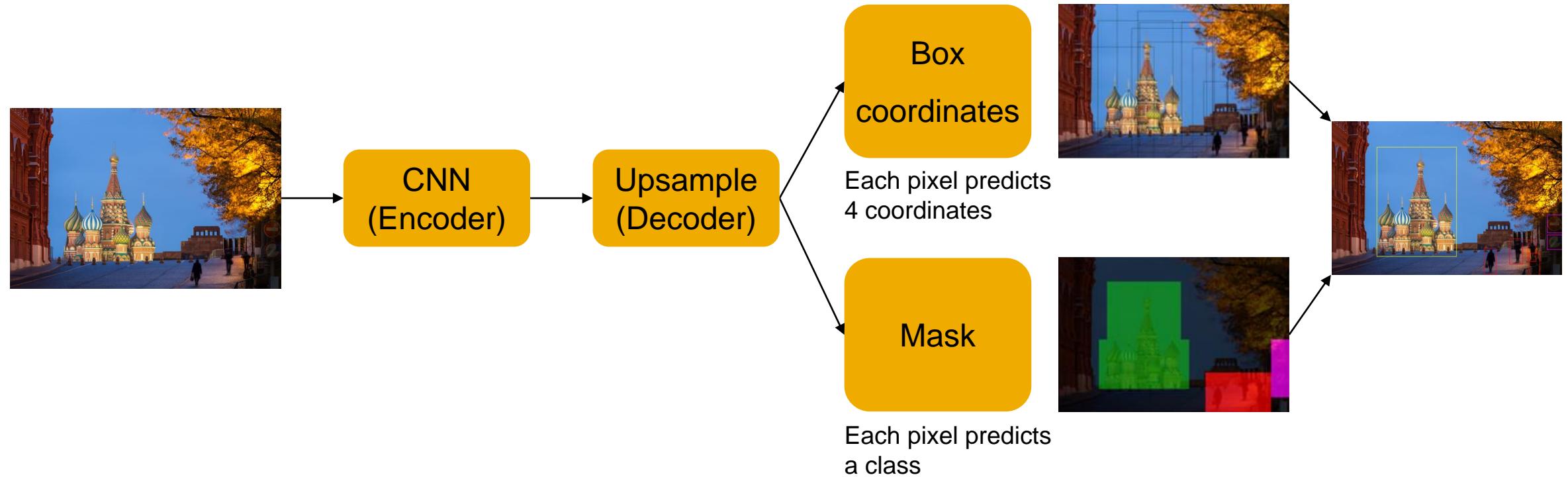
# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]



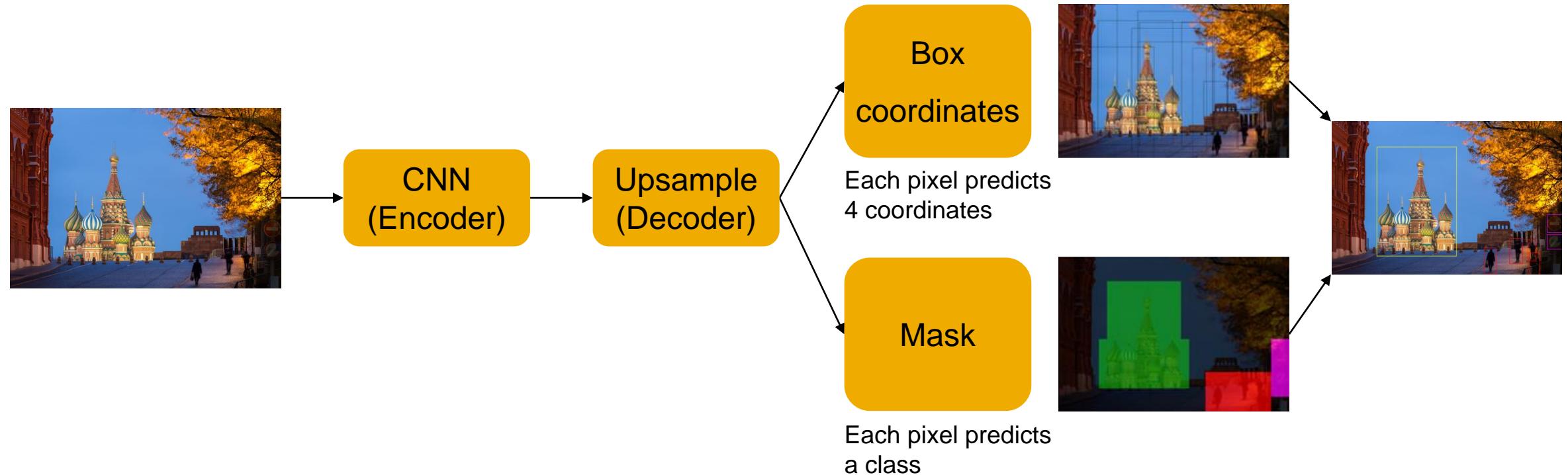
# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]



# Applications of CNNs

One-stage detectors e.g.: YOLO(9000) – You Only Look Once [Redmon, Divvala, Girshick, Farhadi, 2015]



Also see: Single Shot Multibox Detector (SSD)

# Applications of CNNs

One-stage detectors: comparison

**Problems:** One-stage detectors are less accurate than two-stage detectors:

# Applications of CNNs

## One-stage detectors: comparison

**Problems:** One-stage detectors are less accurate than two-stage detectors:

- Number of bounding box predictions is considerably larger (c.f. R-CNN's 2k proposals vs. 100k in most one-stage detectors)

# Applications of CNNs

One-stage detectors: comparison

**Problems:** One-stage detectors are less accurate than two-stage detectors:

- Number of bounding box predictions is considerably larger (c.f. R-CNN's 2k proposals vs. 100k in most one-stage detectors)
- This leads to more class imbalance since most of the bounding boxes **belong to background**

# Applications of CNNs

One-stage detectors: RetinaNet [Lin, Goyal, Girshick, He, Dollar, 2017]

**Problems:** One-stage detectors are less accurate than two-stage detectors:

- Number of bounding box predictions is considerably larger (c.f. R-CNN's 2k proposals vs. 100k in most one-stage detectors)
- This leads to more class imbalance since most of the bounding boxes **belong to background**

**Remedy:** Introduce so-called *Focal Loss*:

# Applications of CNNs

One-stage detectors: RetinaNet [Lin, Goyal, Girshick, He, Dollar, 2017]

**Problems:** One-stage detectors are less accurate than two-stage detectors:

- Number of bounding box predictions is considerably larger (c.f. R-CNN's 2k proposals vs. 100k in most one-stage detectors)
- This leads to more class imbalance since most of the bounding boxes **belong to background**

**Remedy:** Introduce so-called *Focal Loss*:

- Down-weight contribution of examples that are particularly easy to learn (i.e. background patches)

# Applications of CNNs

One-stage detectors: RetinaNet [Lin, Goyal, Girshick, He, Dollar, 2017]

**Problems:** One-stage detectors are less accurate than two-stage detectors:

- Number of bounding box predictions is considerably larger (c.f. R-CNN's 2k proposals vs. 100k in most one-stage detectors)
- This leads to more class imbalance since most of the bounding boxes **belong to background**

**Remedy:** Introduce so-called *Focal Loss*:

- Down-weight contribution of examples that are particularly easy to learn (i.e. background patches)

$$- \quad L(p, y) = \begin{cases} -\log p, & y = +1 \\ -\log(1 - p), & y = -1 \end{cases}$$

$$- \quad F(p, y) = \begin{cases} -(1 - p)^\gamma \log p, & y = +1 \\ -p^\gamma \log(1 - p), & y = -1 \end{cases}$$

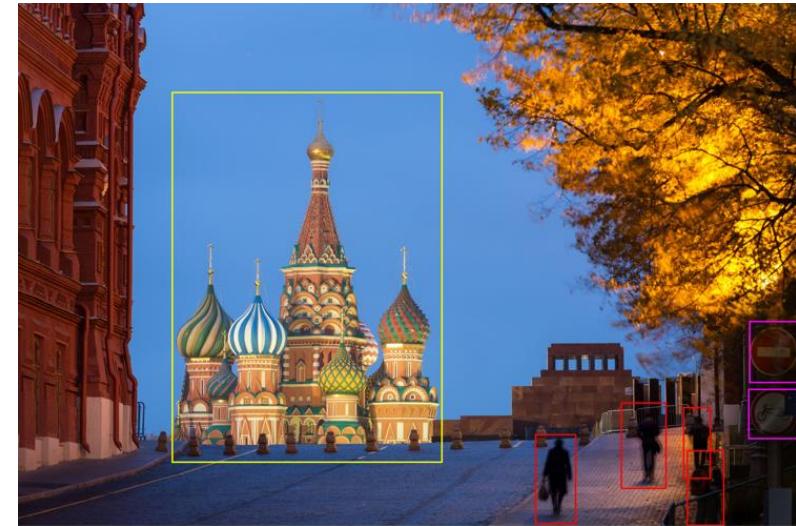
For more details: RetinaNet – <https://arxiv.org/abs/1708.02002>

# Applications of CNNs

## Overview

### Content:

- Object detection
  - Two-stage detectors
  - One-stage detectors
- **Segmentation**



# Applications of CNNs

Image segmentation: example



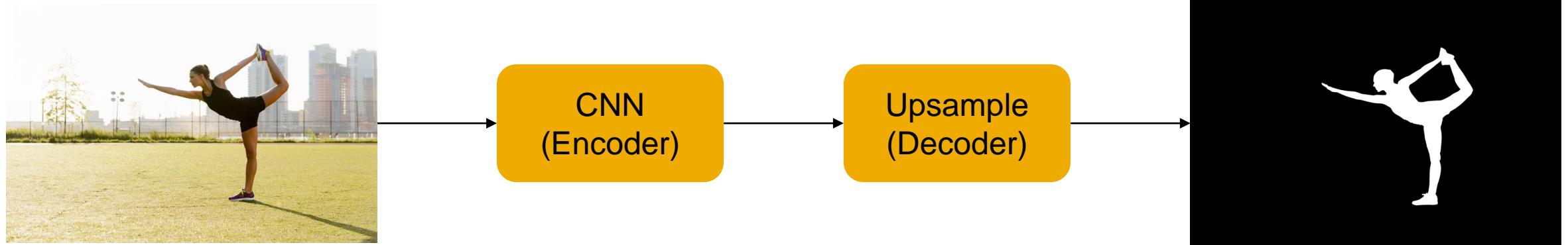
# Applications of CNNs

Image segmentation: example



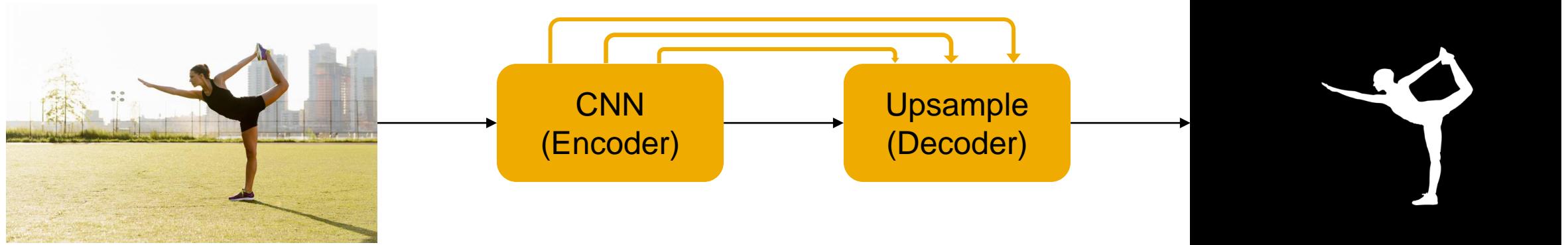
# Applications of CNNs

Image segmentation: generic architecture



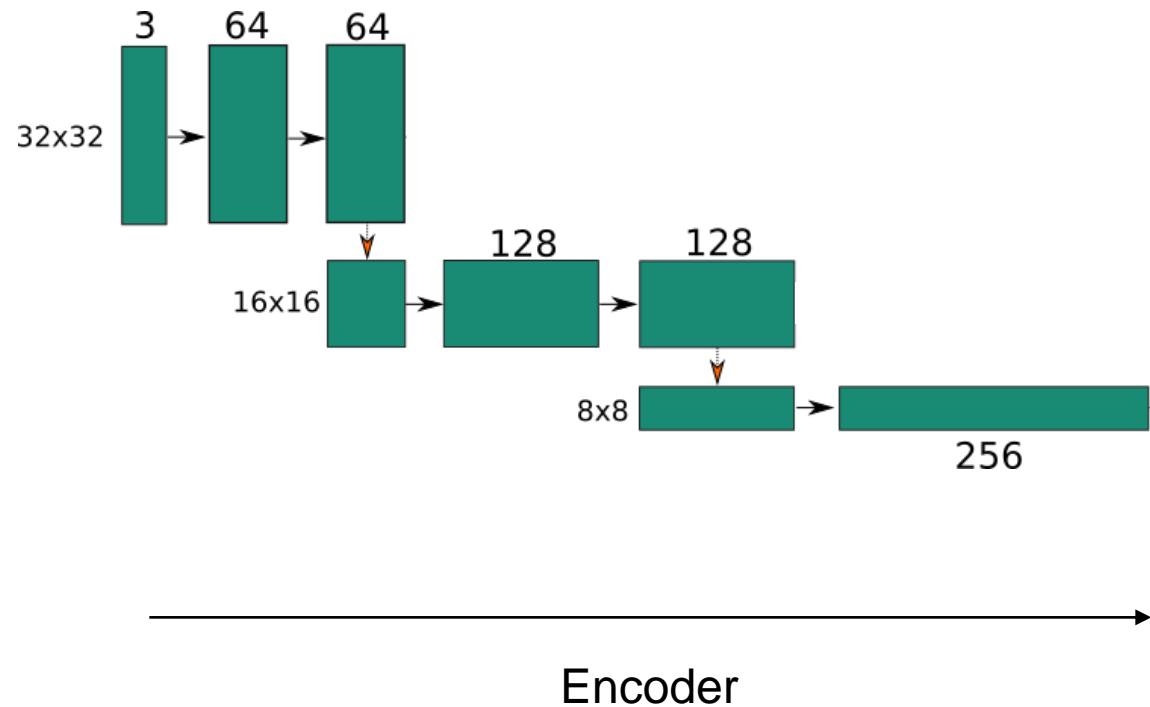
# Applications of CNNs

Image segmentation: generic architecture



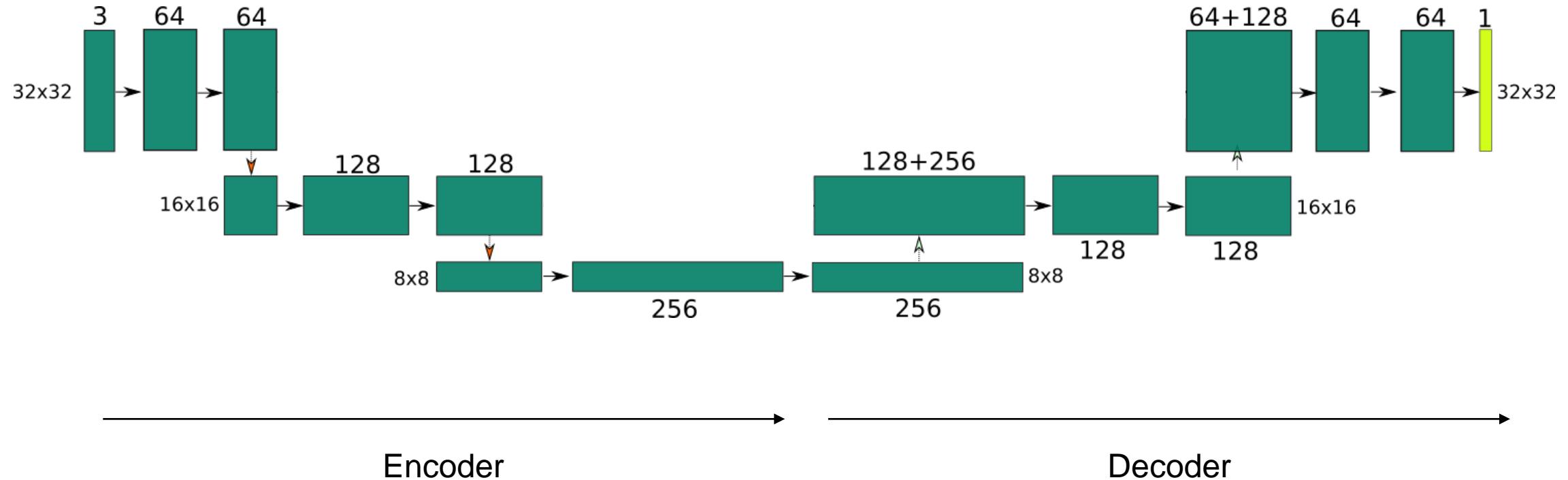
# Applications of CNNs

Image segmentation e.g. U-Net architecture



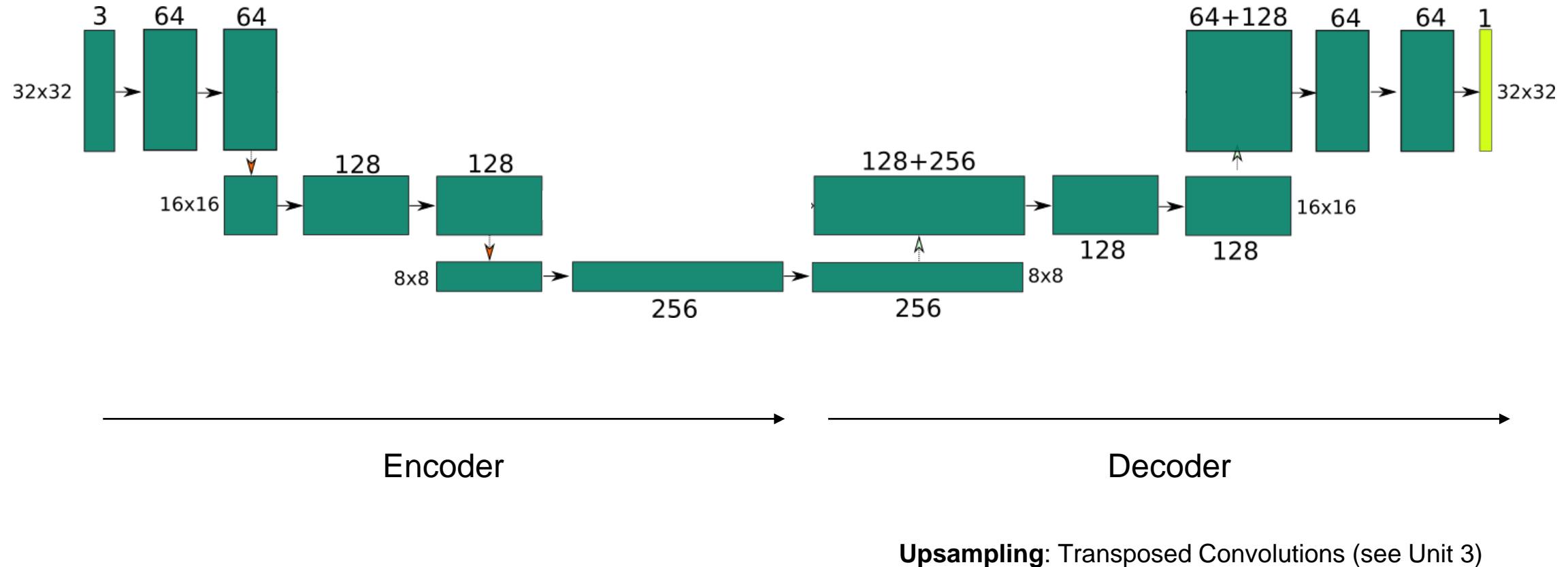
# Applications of CNNs

Image segmentation e.g. U-Net architecture



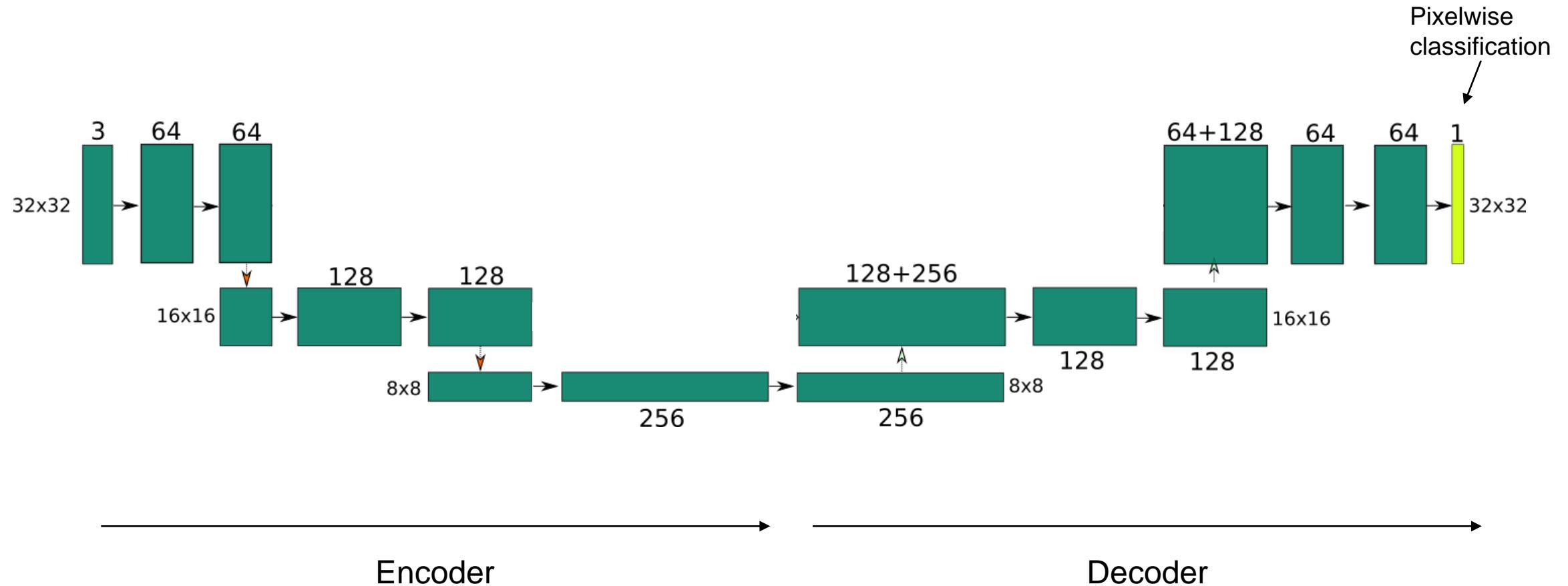
# Applications of CNNs

Image segmentation e.g. U-Net architecture



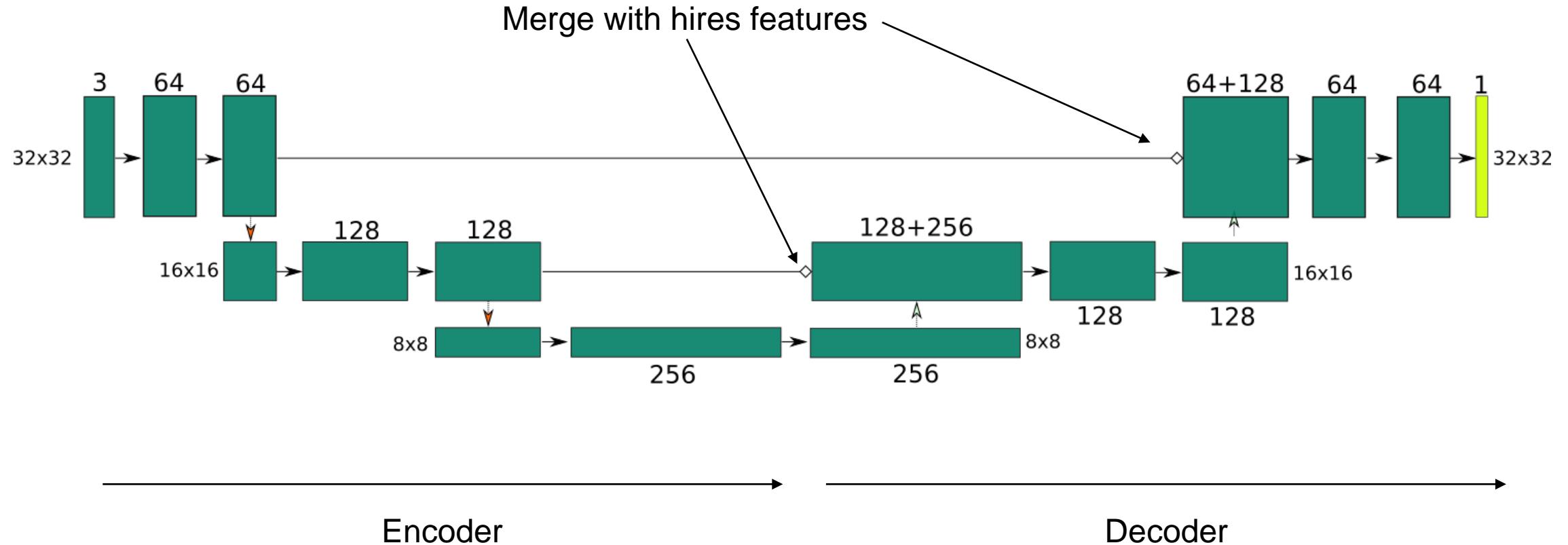
# Applications of CNNs

Image segmentation e.g. U-Net architecture



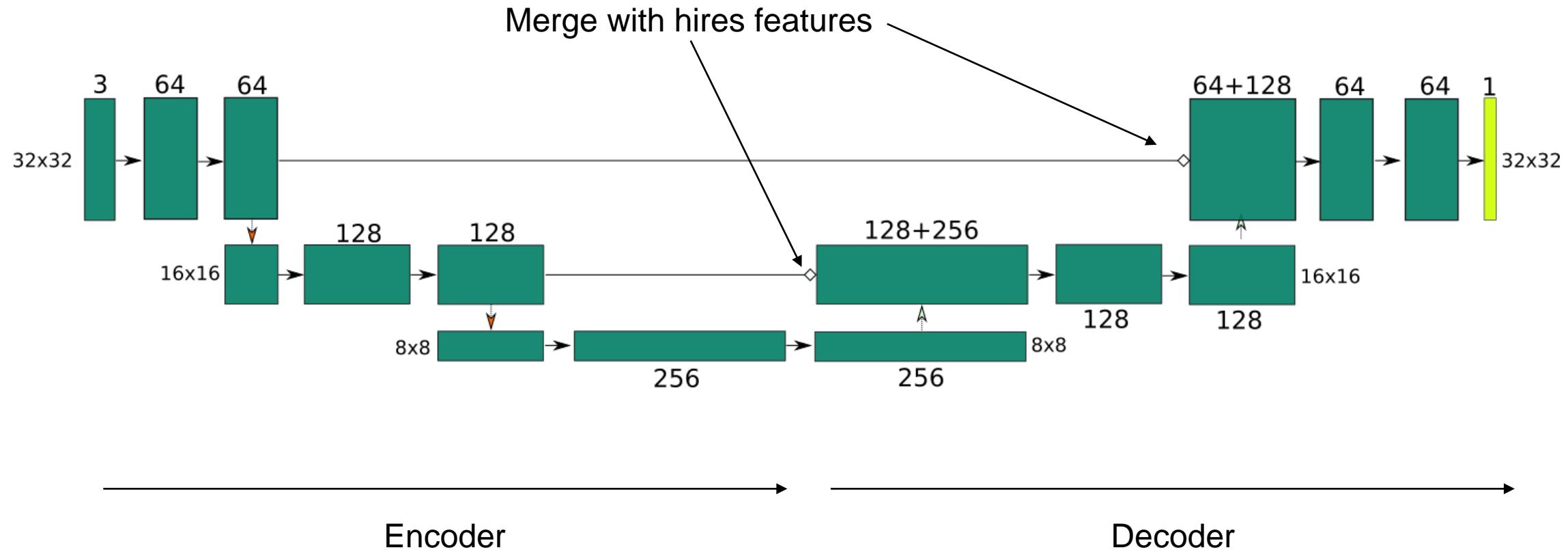
# Applications of CNNs

Image segmentation e.g. U-Net architecture



# Applications of CNNs

Image segmentation e.g. U-Net architecture



Also see: e.g. SegNet  
e.g. Pyramid Scene Parsing Network

# Applications of CNNs

Image segmentation: challenges

## ■ Efficiency / Accuracy

- Encoder downsampling (low resolution) versus dense prediction

# Applications of CNNs

## Image segmentation: challenges

### ■ Efficiency / Accuracy

- Encoder downsampling (low resolution) versus dense prediction
  - Make use of dilated convolutions
  - Feature fusion

# Applications of CNNs

Image segmentation: challenges

## ■ Efficiency / Accuracy

- Encoder downsampling (low resolution) versus dense prediction
  - Make use of dilated convolutions
  - Feature fusion

## ■ Data availability

- Ground truth: Hand-segmented pictures (very costly to acquire!)
-  Data augmentation is essential!

# Applications of CNNs

Image segmentation: data augmentation

- Data augmentation is indispensable when there is little training data
- Use randomized “on-the-fly” modifications

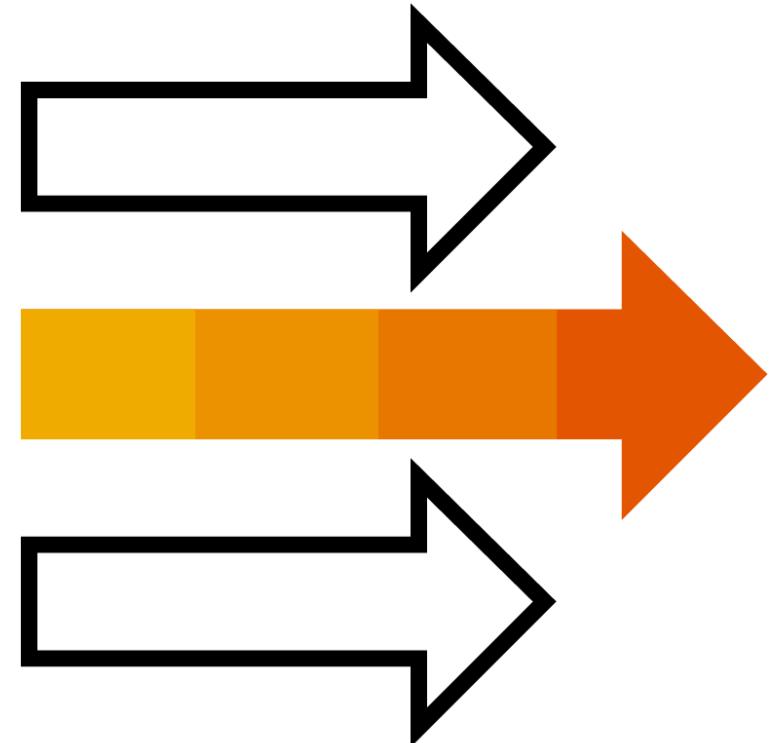


# Applications of CNNs

Coming up next

## Industry Applications of Deep Learning

- Machine learning in customer service
- Machine learning in banking
- Medical image segmentation



# Thank you.

Contact information:

[open@sap.com](mailto:open@sap.com)

# © 2017 SAP SE or an SAP affiliate company. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or for any purpose without the express permission of SAP SE or an SAP affiliate company.

The information contained herein may be changed without prior notice. Some software products marketed by SAP SE and its distributors contain proprietary software components of other software vendors. National product specifications may vary.

These materials are provided by SAP SE or an SAP affiliate company for informational purposes only, without representation or warranty of any kind, and SAP or its affiliated companies shall not be liable for errors or omissions with respect to the materials. The only warranties for SAP or SAP affiliate company products and services are those that are set forth in the express warranty statements accompanying such products and services, if any. Nothing herein should be construed as constituting an additional warranty.

In particular, SAP SE or its affiliated companies have no obligation to pursue any course of business outlined in this document or any related presentation, or to develop or release any functionality mentioned therein. This document, or any related presentation, and SAP SE's or its affiliated companies' strategy and possible future developments, products, and/or platform directions and functionality are all subject to change and may be changed by SAP SE or its affiliated companies at any time for any reason without notice. The information in this document is not a commitment, promise, or legal obligation to deliver any material, code, or functionality. All forward-looking statements are subject to various risks and uncertainties that could cause actual results to differ materially from expectations. Readers are cautioned not to place undue reliance on these forward-looking statements, and they should not be relied upon in making purchasing decisions.

SAP and other SAP products and services mentioned herein as well as their respective logos are trademarks or registered trademarks of SAP SE (or an SAP affiliate company) in Germany and other countries. All other product and service names mentioned are the trademarks of their respective companies.

See <http://global.sap.com/corporate-en/legal/copyright/index.epx> for additional trademark information and notices.